

Balance is Essence: Accelerating Sparse Training via Adaptive Gradient Correction

Bowen Lei^{1*}, Dongkuan Xu², Ruqi Zhang³, Shuren He¹, Bani Mallick¹
¹Texas A&M University, ²North Carolina State University, ³Purdue University
*bowenlei@stat.tamu.edu

Despite impressive performance, deep neural networks require significant memory and computation costs, prohibiting their application in resource-constrained scenarios. Sparse training is one of the most common techniques to reduce these costs, however, the sparsity constraints add difficulty to the optimization, resulting in an increase in training time and instability. In this work, we aim to overcome this problem and achieve space-time co-efficiency. To accelerate and stabilize the convergence of sparse training, we analyze the gradient changes and develop an adaptive gradient correction method. Specifically, we approximate the correlation between the current and previous gradients, which is used to balance the two gradients to obtain a corrected gradient. Our method can be used with the most popular sparse training pipelines under both standard and adversarial setups. Theoretically, we prove that our method can accelerate the convergence rate of sparse training. Extensive experiments on multiple datasets, model architectures, and sparsities demonstrate that our method outperforms leading sparse training methods by up to 5.0% in accuracy given the same number of training epochs, and reduces the number of training epochs by up to 52.1% to achieve the same accuracy. Our code is available on: <https://github.com/StevenBoys/AGENT>.

1. Introduction

Sparse training [1–3] is one of the most popular classes of methods to improve the efficiency of deep neural networks (DNNs) in terms of space (e.g. memory storage), and it is receiving increasing attention, especially in resource-limited situations [2–4]. During sparse training, a certain percentage of connections are removed to save memory [2, 4]. Sparse patterns, which describe where connections are retained or removed, are iteratively updated [2, 5–7]. The goal is to find a resource-efficient sparse neural network (i.e., removing some connections) with comparable or even higher performance compared to the original dense model (i.e., keeping all connections) [8–10].

However, sparse training can bring some side effects to the training process, especially in the case of high sparsity (e.g., 99% weights are zero). First, sparsity can increase the variance of stochastic gradients, leading the model to move in a sub-optimal direction and hence slow convergence [11, 12]. As shown in Figure 1 (a), we empirically see that the gradient variance grows with increasing sparsity (more details in Section C.1). Second, it can result in training instability (i.e., a noisy trajectory of test accuracy w.r.t. iterations) [13, 14], which requires additional time to compensate for the accuracy drop, resulting in slow convergence [15]. Additionally, the need to consider the robustness of the model during sparse training is highlighted in order to apply sparse training to a wide range of real-world scenarios where there are often challenges with dataset shifts [7, 11, 16, 17]. To address these issues, we raise the following questions:

Question 1. How to simultaneously improve convergence speed and training stability of sparse training?

Prior gradient correction methods [18–20] are used to accelerate and stabilize dense training, while we find that it fails in sparse training. They usually assume that current and previous gradients are highly correlated, and therefore add a large constant amount of previous gradients to correct the gradient [19, 21, 22]. However, this assumption does not hold in sparse training. Figure 1 (b) shows

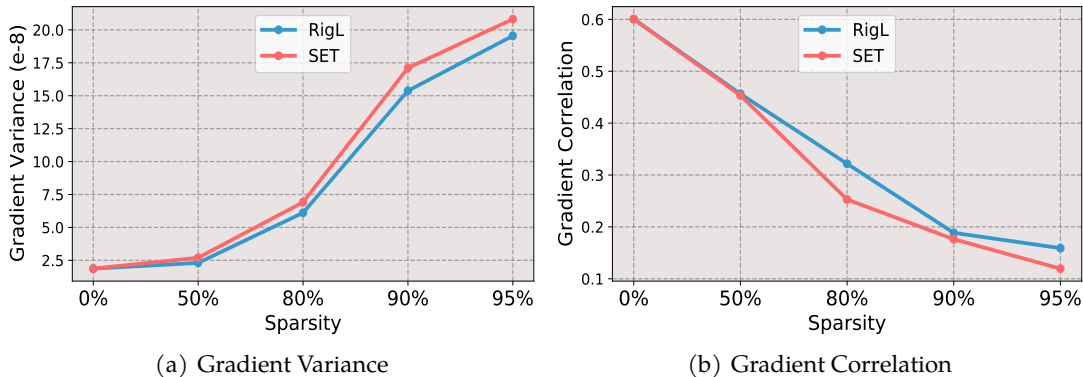


Figure 1: Gradient variance (a) and gradient correlation (b) of models obtained by RigL and SET at different sparsities including 0% (dense), 50%, 80%, 90%, 95%. Gradient variance grows with increasing sparsity. Gradient correlation drops with increasing sparsity. The sparse models have larger gradient variance and smaller gradient correlation compared to dense models. The gradient correlation at different sparsities, implying that the gradient correlation decreases with increasing sparsity (more details in Section C.1), which breaks the balance between current and previous gradients. Therefore, we propose to adaptively change the weights of previous and current gradients based on their correlation to add an appropriate amount of previous gradients.

Question 2. *How to design an accelerated and stabilized sparse training method that is effective in real-world scenarios with dataset shifts?*

Moreover, real-world applications are under-studied in sparse training. Prior methods use adversarial training to improve model robustness and address the challenge of data shifts, which usually introduces additional bias beyond the variance in the gradient estimation [23], increasing the difficulty of gradient correction (more details in Section 4.2). Thus, to more accurately approximate the full gradient, especially during the adversarial setup, we design a scaling strategy to control the weights of the two gradients, determining the amount of previous gradient information to be added to the current gradient, which helps the balance and further accelerates the convergence.

In this work, we propose an **adaptive gradient correction (AGENT)** method to accelerate and stabilize sparse training for both standard and adversarial setups. Theoretically, we prove that our method can accelerate the convergence rate of sparse training. Empirically, we perform extensive experiments on multiple benchmark datasets, model architectures, and sparsities. In both standard and adversarial setups, our method improves the accuracy by up to **5.0%** given the same number of epochs and reduces the number of epochs up to **52.1%** to achieve the same performance compared to the leading sparse training methods. In contrast to previous efforts of sparse training acceleration which mainly focus on structured sparse patterns, our method is compatible with both unstructured and structured sparse training pipelines [24, 25].

2. Related Work

Sparse Training: Interest in sparse DNNs has been on the rise recently. The goal is to achieve comparable performance with sparse weights to satisfy the constraints. Different sparse training methods have emerged, where sparse weights are maintained in the training process. Various pruning and growth criteria are proposed, such as weight/gradient magnitude, random selection, and weight sign [1–7, 26–33]. However, the aforementioned studies focus on improving the performance, while neglecting the side effect of sparse training. Sparsity not only increases gradient variance, thus delaying convergence [11, 12], but also leads to training instability [14]. It is a challenge to achieve both space and time efficiency. Additionally, sparse training can also exacerbate models’ vulnerability to adversarial samples, which is one of the weaknesses of DNNs [7]. When the model encounters intentionally manipulated data, its performances may deteriorate rapidly, leading to increasing security concerns [34, 35]. In this paper, we focus on sparse training. In general, our method can be applied to any SGD-based sparse training pipelines.

Accelerating Training: Studies have been conducted in recent years on achieving time efficiency in DNNs [36, 37], and one popular direction is to obtain a more accurate gradient estimate to update

the model [20], such as variance reduction. In SGD, one uses small batches of data to approach the full gradient. The batch estimator is usually unbiased but can have a large variance and misguide the model, leading to studies on variance reduction [18–20, 38–43]. While adversarial training brings bias in the estimator [23], we need to face the bias-variance tradeoff when doing gradient correction. A shared idea is to balance the gradient noise with a less noisy old gradient [19, 44, 45]. Some other momentum-based methods have a similar strategy of using old information [46, 47]. However, all the above work considers only the acceleration in non-sparse cases.

Acceleration is more challenging in sparse training, and previous research on it has focused on structured sparse training [24, 25, 48]. First, sparse training will induce larger variance [11]. In addition, some key assumptions associated with gradient correction methods do not hold under sparsity constraints. In the non-sparse case, the old and new gradients are assumed to be highly correlated, so we can collect a large amount of knowledge from the old gradients [19, 21, 22]. However, sparsity tends to lead to lower correlations, and this irrelevant information can be harmful, making previous methods no longer applicable to sparse training and requiring a finer balance between new and old gradients. Furthermore, the structured sparsity pattern is not flexible enough, which can lead to lower model accuracy. In contrast, our method accelerates sparse training from an optimization perspective and is compatible with both unstructured and structured sparse training pipelines.

3. Preliminaries: Stochastic Variance Reduced Gradient

Stochastic variance reduced gradient (SVRG) [21, 39, 49] is a widely-used gradient correction method designed to obtain more accurate gradient estimates, which has been followed by many studies [18, 19, 50]. Specifically, after each epoch of training, we evaluate the full gradients $\tilde{\mathbf{g}}$ based on $\tilde{\theta}$ at that time and store them for later use. In the next epoch, the batch gradient estimate on \mathbf{B}_t is updated using the stored old gradients via Eq. (1).

$$\hat{\mathbf{g}}(\theta_t) = \frac{1}{n} \sum_{i \in \mathbf{B}_t} \left(\mathbf{g}_i(\theta_t) - \mathbf{g}_i(\tilde{\theta}) \right) + \tilde{\mathbf{g}} \quad (1)$$

where $\mathbf{g}_i(\theta_t) = \nabla l(\mathbf{x}_i | \theta_t)$, $l(\theta_t) = (\sum_{i=1}^N l(\mathbf{x}_i | \theta_t)) / N$ is the loss function, $\tilde{\mathbf{g}} = \nabla l(\tilde{\theta})$, θ_t is the current parameters, n is the number of samples in each mini-batch data, and N is the total number of samples. SVRG successfully accelerates many training tasks in the non-sparse case, but does not work well in sparse training, which is similar to many other gradient correction methods.

4. Method

We propose an adaptive gradient correction (AGENT) method and integrate it with recent sparse training pipelines to achieve accelerations and improve training stability. To accomplish the goal, our AGENT filters out less relevant information and obtains a well-controlled and time-varying amount of knowledge from the old gradients. Our method overcomes the limitations of previous acceleration methods such as SVRG [21, 49, 51], and successfully accelerates and stabilizes sparse training. Our AGENT method is outlined in Algorithm 1 and illustrated in the following sections.

4.1. Adaptive Control over Old Gradients

In AGENT, we designed an adaptive addition of old gradients to new gradients to filter less relevant information and achieve a balance between new and old gradients. Specifically, we add an adaptive weight $c_t \in [0, 1]$ to the old gradient as shown in Eq. (2), where we use $\mathbf{g}_{\text{new}} = \frac{1}{n} \sum_{i \in \mathbf{B}_t} \mathbf{g}_i(\theta_t)$ and $\mathbf{g}_{\text{old}} = \frac{1}{n} \sum_{i \in \mathbf{B}_t} \mathbf{g}_i(\tilde{\theta})$ to denote the gradient on current parameters θ_t and previous parameters $\tilde{\theta}$ for a random subset \mathbf{B}_t , respectively. When the old and new gradients are highly correlated, we need a large c to get more useful information from the old gradient. Conversely, when the relevance is low, we need a smaller c so that we do not let irrelevant information corrupt the new gradient.

$$\hat{\mathbf{g}}(\theta_t) = \frac{1}{n} \sum_{i \in \mathbf{B}_t} \left(\mathbf{g}_i(\theta_t) - c_t \cdot \mathbf{g}_i(\tilde{\theta}) \right) + c_t \cdot \tilde{\mathbf{g}} = \mathbf{g}_{\text{new}} - c_t \cdot \mathbf{g}_{\text{old}} + c_t \cdot \tilde{\mathbf{g}}. \quad (2)$$

A suitable c_t should effectively reduce the variance of $\hat{\mathbf{g}}(\boldsymbol{\theta}_t)$. We decompose the variance of $\hat{\mathbf{g}}(\boldsymbol{\theta}_t)$ in Eq. (3) with some abuse of notation, where the variance of the updated gradient is a quadratic function of c_t . For simplicity, considering the case where $\hat{\mathbf{g}}(\boldsymbol{\theta}_t)$ is a scalar, the optimal c_t^* will be in the form of Eq. (3). As we can see, c_t^* is not close to 1 when the new gradient is not highly correlated with the old gradient. Since low correlation between \mathbf{g}_{new} and \mathbf{g}_{old} is more common in sparse training, directly setting $c_t = 1$ in previous methods is not appropriate and we need to estimate adaptive weights c_t^* . In support of this claim, we include a discussion and empirical analysis in the Appendix B.6 to demonstrate that as sparsity increases, the gradient changes faster, leading to lower correlations between \mathbf{g}_{new} and \mathbf{g}_{old} .

$$\text{Var}(\hat{\mathbf{g}}(\boldsymbol{\theta}_t)) = \text{Var}(\mathbf{g}_{\text{new}}) + c_t^2 \cdot \text{Var}(\mathbf{g}_{\text{old}}) - 2c_t \cdot \text{Cov}(\mathbf{g}_{\text{new}}, \mathbf{g}_{\text{old}}), \quad c_t^* = \frac{\text{Cov}(\mathbf{g}_{\text{new}}, \mathbf{g}_{\text{old}})}{\text{Var}(\mathbf{g}_{\text{old}})}. \quad (3)$$

We find it impractical to compute the exact c_t^* and thus propose an approximation algorithm for it to obtain a balance between the new and old gradient. There are two challenges to calculate the exact c_t^* . On the one hand, to approach the exact value, we need to calculate the gradients on every batch data, which is too expensive to do in each iteration. On the other hand, the gradients are often high-dimensional and the exact optimal c_t^* will be different for different gradients. Thus, inspired by Deng et al. [52], we design an approximation algorithm that makes good use of the loss information and leads to only a small increase in computational effort. More specifically, we estimate c_t^* according to the changes of loss as shown in Eq. (4) and update \hat{c}_t^* adaptively before each epoch using momentum. Loss is a scalar, which makes it possible to estimate the shared correlation for all current and previous gradients. In addition, the loss is intuitively related to gradients and the correlation between losses can give us some insights into that of the gradients.

$$\hat{c}_t^* = \frac{\text{Cov}(l(\mathbf{B}|\boldsymbol{\theta}_t), l(\mathbf{B}|\tilde{\boldsymbol{\theta}}))}{\text{Var}(l(\mathbf{B}|\tilde{\boldsymbol{\theta}}))}, \quad (4)$$

where \mathbf{B} denotes a subset of samples used to estimate the gradients.

4.2. Additional Scaling Parameter is Important

To guarantee successful acceleration in sparse and adversarial training, we further propose a scaling strategy that multiplies the estimated c_t^* by a small scaling parameter γ . There are two main benefits of using a scaling parameter. First, the scaling parameter γ can reduce the bias of the gradient estimates in adversarial training [23]. In standard training, the batch gradient estimator is an unbiased estimator of the full gradient. However, in adversarial training, we perturb the mini-batch of samples \mathbf{B}_t into $\tilde{\mathbf{B}}_t$. The old gradients \mathbf{g}_{old} are calculated on batch data \mathbf{B}_t , but the stored old gradients $\tilde{\mathbf{g}}$ are obtained from the original data including $\tilde{\mathbf{B}}_t$, which makes $\mathbb{E}[\mathbf{g}_{\text{old}} - \tilde{\mathbf{g}}]$ unequal to zero. Consequently, as shown in Eq. (5), the corrected estimator for full gradients will no longer be unbiased. It may have a small variance but a large bias, resulting in poor performance. Therefore, we propose a scaling parameter γ between 0 and 1 to reduce the bias from $c_t(\mathbf{g}_{\text{old}} - \tilde{\mathbf{g}})$ to $\gamma c_t(\mathbf{g}_{\text{old}} - \tilde{\mathbf{g}})$.

$$\mathbb{E}[\hat{\mathbf{g}}(\boldsymbol{\theta}_t)] = \mathbb{E}[\mathbf{g}_{\text{new}} - c_t(\mathbf{g}_{\text{old}} - \tilde{\mathbf{g}})] \neq \mathbb{E}[\mathbf{g}_{\text{old}}] = \frac{1}{N} \sum_{i=1}^N \mathbf{g}_i(\boldsymbol{\theta}_t). \quad (5)$$

Second, the scaling parameter γ guarantees that the variance can still be reduced in the face of worst-case estimates of c_t^* to accelerate the training. The key idea is illustrated in Figure 2, where

Algorithm 1 Adaptive Gradient Correction

Input: $\tilde{\boldsymbol{\theta}} = \boldsymbol{\theta}_0$, epoch length m , step size η_t , $c_0 = 0$, scaling parameter γ , smoothing factor α
for $t = 0$ **to** $T - 1$ **do**
 if $t \bmod m = 0$ **then**
 $\tilde{\boldsymbol{\theta}} = \boldsymbol{\theta}_t$
 $\tilde{\mathbf{g}} = (\sum_{i=1}^N \nabla l(\mathbf{x}_i|\tilde{\boldsymbol{\theta}}))/N$
 if $t > 0$ **then**
 Calculate \hat{c}_t^* via Eq. (4)
 $c_t = (1 - \alpha)c_{t-1} + \alpha\hat{c}_t^*$
 end if
 else
 $c_t = c_{t-1}$
 end if
 Sample a mini-batch data \mathbf{B}_t with size n
 $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta_t \cdot \left(\frac{1}{n} \sum_{i \in \mathbf{B}_t} (\mathbf{g}_i(\boldsymbol{\theta}_t) - \gamma c_t \cdot \mathbf{g}_i(\tilde{\boldsymbol{\theta}})) + \gamma c_t \cdot \tilde{\mathbf{g}} \right)$
end for

x and y axis correspond to the weight c_t and the gradient variance, respectively. The blue curve is a quadratic function that represents the relationship between c_t and the variance. Suppose the true optimal is c^* , and we make an approximation to it. In the worst case, this approximation may be as bad as \hat{c}_1 , making the variance even larger than a_3 (variance in SGD) and slowing down the training. Then, if we replace \hat{c}_1 with $\gamma\hat{c}_1$, we can reduce the variance and accelerate the training.

4.3. Connection to Other Optimizers

Momentum-based Methods: Our AGENT is designed with a similar idea to the momentum-based method [53, 54], where old gradients are used to improve the current batch gradient. However, the momentum-based method does not consider sparse and adversarial training characteristics such as the reduced correlation between current and previous gradients and potential bias of gradient estimator, and fails to provide an adaptive balance between old and new information. When the correlation is low, it can still incorporate too much of the old information and increase the gradient variance or bias.

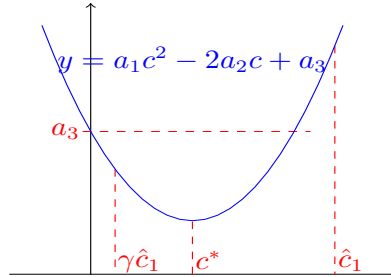


Figure 2: Illustration of how the scaling parameter $\gamma = 0.1$ ensures the acceleration in the face of worst-case estimate of c_t^* . The blue curve is a quadratic function, representing the relationship between c_t and the variance. c^* is the optimal value. \hat{c}_1 is a poor estimate making the variance larger than a_3 (variance in SGD). $\gamma\hat{c}_1$ can reduce the variance.

Adaptive Gradient Method: Our AGENT can be viewed as a new type of adaptive gradient method that adaptively adjusts the amount of gradient information used to update parameters, such as Adam [55]. However, previous methods are not designed for sparse training. Despite their adaptive gradients, their adaptivity is different and does not take the reduced correlation into account.

On the contrary, our AGENT is tailored to the characteristics of sparse training, which approximates the correlation and adds an adaptive weight to the old gradient to establish a balance between the old and new gradients (see more comparisons in Appendix B.9, B.10).

5. Theoretical Justification

Theoretically, we provide a convergence analysis for our AGENT and compare it to SVRG [56]. We use $l(\cdot)$ to denote the loss function and g to denote the gradient. Our proof is based on Assumptions 1-2, and detailed derivation is included in Appendix A.

Assumption 1. (*L-smooth*): The differentiable loss function $l : \mathbb{R}^n \rightarrow \mathbb{R}$ is *L-smooth*, i.e., for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$, the loss l satisfies $\|\nabla l(\mathbf{x}) - \nabla l(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|$.

Assumption 2. (*σ -bounded*): The loss function l has a σ -bounded gradient, i.e., $\|\nabla l_i(\mathbf{x})\| \leq \sigma$ for all $i \in [N]$ and $\mathbf{x} \in \mathbb{R}^n$.

Given Assumptions 1-2, we follow the analysis framework above and establish Theorem 1.

Theorem 1. Under Assumptions 1-2, with proper choice of step size η_t and c_t , the gradient $\mathbb{E}[\|g(\boldsymbol{\theta}_\pi)\|^2]$ using AGENT after T training epochs can be bounded by:

$$\mathbb{E}[\|g(\boldsymbol{\theta}_\pi)\|^2] \leq \frac{(l(\boldsymbol{\theta}_0) - l(\boldsymbol{\theta}_*))LN^\alpha}{Tn\nu} + \frac{2\kappa\mu^2\sigma^2}{N^\alpha m\nu}$$

where $\boldsymbol{\theta}_\pi$ is sampled uniformly from $\{\{\boldsymbol{\theta}_t^s\}_{t=0}^{m-1}\}_{s=0}^{T-1}$, N denotes the data size, n denotes the mini-batch size, m denotes the epoch length, $\boldsymbol{\theta}_0$ is the initial point and $\boldsymbol{\theta}_*$ is the optimal solution, $\nu, \mu, \kappa, \alpha > 0$ are constants depending on η_t and c_t , N and n .

In regard to Theorem 1, we make the following remarks to justify the acceleration from our AGENT:

Remark 1. (*Faster Gradient Change Speed*) An influential difference between sparse and dense training is the gradient change speed, which is reflected in Assumption 1 (*L-smooth*). Typically, L in sparse training will be larger than L in dense training.

Remark 2. (First Term Analysis) In Theorem 1, the first term in the bound of our AGENT measures the error from deviations of the optimal parameters, which goes to zero when the number of epochs T reaches infinity. However, in real sparse training applications, T is finite and this term is expanded due to the increase of L in sparse training, which implies that the optimization under sparse constraints is more challenging.

Remark 3. (Second Term Analysis) In Theorem 1, the second term measures the error from the noisy gradient and the finite data in optimization. Since σ^2 is relatively small and N is usually large in our DNNs training, the second term is negligible or much smaller compared to the first term when T is assumed to be finite.

From the above analysis, we can compare the bounds of AGENT and SVRG and find that in the case of sparse training, an appropriate choice of c_t can make the bound for our AGENT tighter than the bound for SVRG by well-corrected gradients.

Remark 4. (Comparison with SVRG) Under Assumptions 1-2, the gradient $\mathbb{E}[\|g(\theta_\pi)\|^2]$ using SVRG after T training epochs can be bounded by [56]:

$$\mathbb{E}[\|g(\theta_\pi)\|^2] \leq \frac{(l(\theta_0) - l(\theta_*))LN^\alpha}{T\nu^*}.$$

This bound is of a similar form to the first term in Theorem 1. Since the second term of Theorem 1 is negligible, we only need to compare the first term. With a proper choice of c_t , the variance of $\hat{g}(\theta_t)$ will decrease, which leads to a smaller ν for AGENT than ν^* for SVRG (details in Appendix A Remark 6). Thus, AGENT can bring a smaller first term compared to SVRG, which indicates a tighter bound of AGENT compared to SVRG.

6. Experiments

We add our AGENT to four recent sparse training pipelines, namely SET [1], RigL [2], BSR-Net [7] and ITOP [6] (see description in Appendix D). Detailed information about the dataset, model architectures, and other training and evaluation setups is provided below.

Datasets & Model Architectures: For datasets, we use CIFAR-10, CIFAR-100 [57], SVHN [58], and ImageNet-2012 [59]. For model architectures, we use VGG-16 [60], ResNet-18, ResNet-50 [61], and Wide-ResNet-28-4 [62].

Training Settings: For sparse training, we choose two sparsity levels (90% & 99%). For BSR-Net, we consider both standard and adversarial setups. In SET, RigL, and ITOP, we focus on standard training. In standard training, we only use the original data instead of using perturbed samples. For adversarial part, we use the perturbed data with two popular objectives (AT and TRADES) [63, 64] and follow the evaluation in BSR-Net [7].

Table 1: Testing accuracy (%) of BSR-Net-based models. Sparse VGG-16 is learned in standard setups. For the same training epochs, Ours often has higher accuracy compared to BSR-Net.

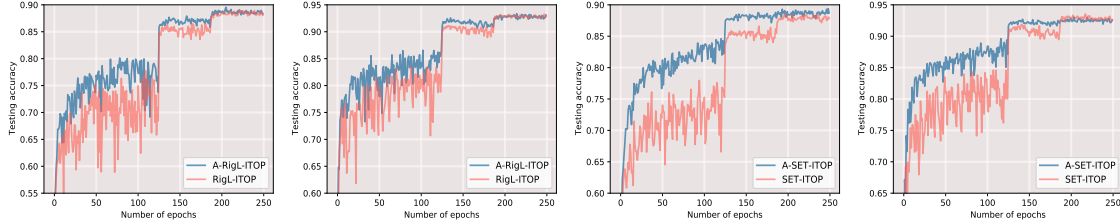
	Epoch	90% Sparsity		99% Sparsity	
		BSR-Net	Ours	BSR-Net	Ours
AT	20	55.0 (1.59)	63.6 (1.31)	49.8 (1.46)	56.4 (1.39)
	40	62.2 (1.88)	64.9 (0.81)	54.1 (1.72)	57.7 (0.39)
	70	73.1 (0.39)	75.1 (0.27)	64.7 (0.30)	66.0 (0.23)
	90	73.2 (0.29)	74.1 (0.25)	63.7 (0.25)	65.8 (0.24)
	140	76.7 (0.27)	77.4 (0.26)	68.4 (0.20)	69.8 (0.14)
	200	76.6 (0.25)	78.1 (0.24)	69.0 (0.15)	70.7 (0.06)
TRADES	20	62.0 (0.82)	65.0 (0.61)	55.7 (0.76)	57.6 (0.45)
	40	65.4 (0.97)	66.0 (0.34)	60.6 (0.69)	58.4 (0.34)
	70	73.4 (0.52)	73.5 (0.33)	66.3 (0.35)	67.3 (0.30)
	90	73.0 (0.36)	73.6 (0.28)	66.2 (0.33)	67.5 (0.24)
	140	76.4 (0.25)	76.8 (0.25)	70.0 (0.29)	69.9 (0.21)
	200	75.6 (0.23)	77.0 (0.24)	70.8 (0.19)	70.9 (0.25)
Standard	20	70.4 (2.50)	81.8 (0.62)	60.6 (1.26)	69.8 (1.45)
	40	77.6 (1.39)	82.4 (0.47)	62.6 (2.47)	73.7 (0.36)
	70	86.8 (0.78)	89.7 (0.38)	79.7 (0.72)	83.7 (0.24)
	90	87.6 (0.63)	89.3 (0.22)	80.5 (0.55)	83.9 (0.42)
	140	91.7 (0.44)	92.5 (0.06)	85.7 (0.42)	86.9 (0.07)
	200	91.8 (0.23)	92.6 (0.12)	85.8 (0.12)	87.1 (0.25)

6.1. Convergence Speed & Stability Comparisons

We compare the convergence speed by two criteria, including (a) **the test accuracy** at the same number of pass data (epoch) and (b) **the number of pass data (epoch)** required to achieve the same test accuracy, which is widely used to compare the speed of optimizers [18, 22, 46, 49].

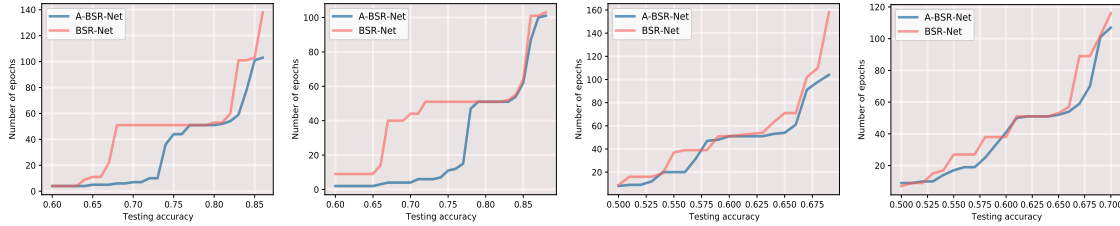
Table 2: Testing accuracy (%) of BSR-Net-based models (BSR) on adversarial samples with VGG-16. Given the training epochs, Ours has higher accuracy compared to BSR-Net in almost all cases.

Epoch	90% Sparse		99% Sparse		90% Sparse		99% Sparse		
	BSR	Ours	BSR	Ours	BSR	Ours	BSR	Ours	
70	37.8	45.2	34.9	39.4	TRADES	34.8	45.4	33.5	39.0
90	33.6	44.8	35.8	39.8		36.8	44.8	31.7	39.1
140	46.5	43.8	40.8	41.2		45.1	46.3	38.2	41.5
200	43.3	44.6	42.2	42.0		47.2	46.2	39.3	41.2



(a) VGG-C(RigL) (b) ResNet-34(RigL) (c) VGG-C(SET) (d) ResNet-34(SET)

Figure 3: Testing accuracy for ITOP-based models at 99% sparsity on CIFAR-10. A-RigL-ITOP and A-SET-ITOP (blue curves) converge faster than RigL-ITOP and SET-ITOP (pink curves).



(a) VGG-16, Standard (b) WRN-28-4, Standard (c) VGG-16, AT (d) WRN-28-4, AT

Figure 4: Number of training epochs required to achieve the accuracy at 99% sparsity. Our A-BSR-Net (blue curves) needs less time to achieve accuracy compared to BSR-Net (pink curves).

Test accuracy at the same number of pass data (epoch): For BSR-Net-based results, Tables 1-2 list the accuracies on clean and adversarial samples of CIFAR-10, for sparse VGG-16, where the higher accuracies are bolded. In the standard setup, we only present clean accuracy. Our method maintains higher clean and robust accuracies for almost all training epochs and setups demonstrating the successful acceleration from our method. In particular, for limited time periods like 20 epochs, our A-BSR-Net usually shows dramatic improvements with clean accuracy as high as **11.4%**, indicating a significant reduction in early search time. In addition, considering the average accuracy improvement over the 6 time budgets, our method outperforms BSR-Net in accuracy by up to **5.0%**.

For ITOP-based results, as shown in Figure 3, the blue curves (A-RigL-ITOP and A-SET-ITOP) are always higher than the pink curves (RigL-ITOP and SET-ITOP), indicating faster training when using our AGENT. In addition, we can see that the pink curves experience severe up-and-down fluctuations, especially in the early stages of training. In contrast, the blue curves are more stable in all the settings, which indicates AGENT is effective in stabilizing the sparse training.

The number of pass data (epoch) required to achieve the same test accuracy: Figure 4 depicts the number of training epochs required to achieve certain accuracy. The blue curves (A-BSR-Net) are always lower than the pink curves (BSR-Net), and on average our method reduces the number of training epochs by up to **52.1%**, indicating faster training when using our proposed A-BSR-Net.

6.2. Final Accuracy Comparisons

In addition, we compare the final accuracy after sufficient training. For ITOP-based results in Table 3, we compare our A-RigL-ITOP with RigL-ITOP on ImageNet-12 using ResNet-50, and ours always maintain the final accuracy. For BSR-Net-based results in Table 4, we compare our A-BSR-Net

Table 3: Final accuracy (%) of ITOP-based ResNet-50 at on ImageNet-2012. Ours maintains the accuracy.

SPARSITY	RigL-ITOP	Ours
80%	75.5 (0.10)	75.6 (0.12)
90%	73.6 (0.12)	73.4 (0.11)

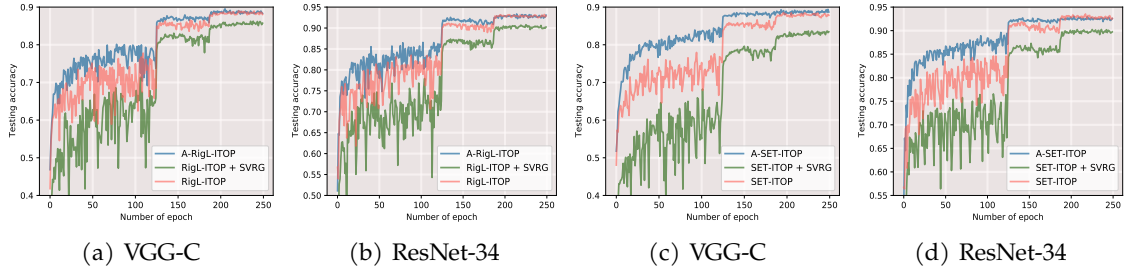


Figure 5: Testing accuracy for ITOP-based models (99%, CIFAR-10). Compared to SGD (pink curves), while SVRG (green curves) slows down the training, ours (blue curves) accelerates it. with BSR-Net on SVHN using VGG-16 and WRN-28-4, and our method is often the best. This shows that our AGENT accelerates sparse training while maintaining or even improving accuracy.

6.3. Comparison with Other Gradient Correction Methods

We also compare our AGENT with SVRG [50], a popular gradient correction method in the non-sparse case. The presented ITOP-based results are based on sparse (99%) VGG-C and ResNet-34 on CIFAR-10. Fig-

ures 5 (a)-(b) show the testing accuracy of A-RigL-ITOP (blue), RigL-ITOP (pink), and RigL-ITOP+SVRG (green). We can see that the green curve for RigL-ITOP+SVRG is often lower than the other two curves for A-RigL-ITOP and RigL-ITOP, indicating that model convergence is slowed down by SVRG. As for the blue curve for our A-RigL-ITOP, it is always on the top of the pink curve for RigL-ITOP and also smoother than the green curve for RigL-ITOP+SVRG, indicating a successful acceleration and stabilization. The SET-ITOP-based results depicted in Figure 5 (c)-(d) show a similar pattern. The green curve (SET-ITOP+SVRG) is often lower than the blue (A-SET-ITOP) and pink (SET-ITOP) curves. This demonstrates that SVRG does not work for sparse training, while our AGENT overcomes its limitations, leading to accelerated and stabilized sparse training.

Table 4: Final accuracy (%) of BSR-Net-based models at 90% and 99% sparsity on SVHN with adversarial objectives (TRADES). Our AGENT maintains or even improves the accuracy.

	BSR-NET (90%)	Ours (90%)	BSR-NET (99%)	Ours (99%)
VGG-16	89.4 (0.29)	94.4 (0.25)	86.4 (0.25)	90.9 (0.26)
WRN-28-4	92.8 (0.24)	95.5 (0.23)	89.5 (0.22)	92.2 (0.19)

6.4. Combination with Other Gradient Correction Methods

In addition to working with SVRG, our AGENT can be combined with other gradient correction methods to achieve sparse training acceleration, such as the momentum-based variance reduction method (MVR) [46].

We train CIFAR-10 on 99% SET-ITOP-based sparse VGG-C using MVR and MVR+AGENT, respectively. As shown in Table 5, MVR+AGENT usually achieves higher test accuracy than MVR for different epochs, which demonstrates the acceleration effect and the generality of our AGENT.

Table 5: Testing accuracy comparisons between MVR and AGENT+MVR. AGENT accelerates MVR in sparse training.

	20-TH	40-TH	70-TH	90-TH	140-TH	200-TH
MVR	62.6	66.8	69.8	71.2	73.5	74.4
AGENT+MVR	71.6	75.7	77.9	79.1	82.3	82.3

6.5. Ablation Studies

We demonstrate the importance of each component in our method AGENT by removing them one by one and comparing the results. Specifically, we consider examining the contribution of the time-varying weight c_t of the old gradients and the scaling parameter γ . The term "Fixed c_t " corresponds to fixing weight $c_t = 0.1$ during training, and "No γ " represents a direct use of \hat{c}_t^* in Eq. (4) and the momentum scheme without adding the scaling parameter γ .

Table 6 shows the clean and robust accuracies of standard and adversarial (AT) training at 90% and 99% sparsity on CIFAR-10 using VGG-16 under different number of training epoch budgets. In the adversarial training (AT and TRADES), we can see that "No γ " is poorly learned and has the worst results. Our method outperforms "Fix c_t " and "No γ " in almost all cases, especially in highly sparse tasks (i.e., 99% sparsity). For standard training, "No γ " can learn some information, but still

Table 6: Ablation Studies: testing accuracy (%) comparisons with Fixed c and No γ on sparse VGG-16. Results are presented as clean/robust accuracy (%). For the same number of training epochs, our method has higher accuracy compared to Fixed c and No γ in almost all cases.

		90% SPARSITY			99% SPARSITY		
		FIXED c_t	NO γ	OURS	FIXED c_t	NO γ	OURS
AT	20-TH	54.1/36.2	28.6/20.1	63.6/37.3	10.0/10.0	10.0/10.0	56.4/31.4
	40-TH	58.9/37.1	20.4/13.0	64.9/37.9	10.0/10.0	10.0/10.0	57.7/34.5
	70-TH	66.8/41.6	19.9/14.7	75.1/45.2	10.0/10.0	10.0/10.0	66.0/39.4
	90-TH	67.7/43.3	21.8/15.6	74.1/44.8	10.0/10.0	10.0/10.0	65.8/39.8
	140-TH	71.4/43.4	20.0/12.1	77.4/43.8	10.0/10.0	10.0/10.0	69.8/41.2
	200-TH	71.7/43.0	20.5/9.5	78.1/44.6	10.0/10.0	10.0/10.0	70.7/42.0
STANDARD	20-TH	80.9/0.0	70.6/0.0	81.8/0.0	73.7/0.0	51.8/0.0	69.8/0.0
	40-TH	83.3/0.0	68.0/0.0	82.4/0.0	74.9/0.0	55.2/0.0	73.7/0.0
	70-TH	90.2/0.0	77.3/0.0	89.7/0.0	84.1/0.0	65.9/0.0	83.7/0.0
	90-TH	89.8/0.0	77.8/0.0	89.3/0.0	80.5/0.0	67.8/0.0	83.9/0.0
	140-TH	92.4/0.0	80.7/0.0	92.5/0.0	87.2/0.0	71.9/0.0	86.9/0.0
	200-TH	92.1/0.0	78.6/0.0	92.6/0.0	86.4/0.0	70.0/0.0	87.1/0.0

performs worse than the other two methods. For "Fix c_t ", it provides a similar convergence speed as our method, while ours tends to have a better final score. Therefore, both the adaptive update of c_t and the multiplication of the scaling parameter γ are important for the acceleration.

6.6. Scaling Parameter Setting

The scaling parameter γ is to avoid introducing large variance due to error in approximating c_t^* and bias due to the adversarial training. The choice of γ is important and can be seen as a hyperparameter tuning process. Our results are based on $\gamma = 0.1$. We check different values from 0 to 1 and find that it is generally better not to set γ to close to 1 or 0. If setting γ close to 1, we will not be able to completely avoid the increase in variance, which leads to a performance drop, similar to "No γ " in Table 6. If γ is set too small, such as 0.01, the weight of the old gradients will be too small and the old gradients will have limited influence on the model update, which will return to SGD's slowdown and training instability (more discussion in Appendix B.3).

7. Discussion and Conclusion

We develop an adaptive gradient correction (AGENT) method for sparse training to achieve time efficiency and reduce training instability from an optimization perspective, which can be incorporated into any SGD-based sparse training pipeline and work in both standard and adversarial setups. To achieve a fine-grained control over the balance of current and previous gradients, we use loss information to analyze gradient changes, and add an adaptive weight to the old gradients. In addition, we design a scaling parameter to reduce the bias of the gradient estimator introduced by the adversarial samples and improve the worst case of the adaptive weight estimate. In theory, we show that our AGENT can accelerate the convergence rate of sparse training. Experiment results on multiple datasets, model architectures, and sparsities demonstrate that our method outperforms state-of-the-art sparse training methods in terms of accuracy by up to **5.0%** and reduces the number of training epochs by up to **52.1%** for the same accuracy achieved.

A number of methods can be employed to reduce the FLOPs in our AGENT. Similar to SVRG, our AGENT increases the training FLOPs in each iteration due to the extra forward and backward used to compute the old gradients. To reduce the FLOPs, the first method is to use sparse gradients [51], which effectively reduces the cost of backward in sparse training and can be easily applied to our method. The second method is parallel computing Allen-Zhu and Hazan [49]. Since the additional forward and backward over the old model parameters are fully parallelizable, we can view it as doubling the mini-batch size. Third, we can follow the idea of SAGA [40] by storing gradients for each single sample. In this way, we do not need extra forward and backward steps, saving the computation. However, it requires extra memory to store the gradients.

References

- [1] Decebal Constantin Mocanu, Elena Mocanu, Peter Stone, Phuong H Nguyen, Madeleine Gibescu, and Antonio Liotta. Scalable training of artificial neural networks with adaptive sparse connectivity inspired by network science. *Nature communications*, 9(1):1–12, 2018.
- [2] Utku Evci, Trevor Gale, Jacob Menick, Pablo Samuel Castro, and Erich Elsen. Rigging the lottery: Making all tickets winners. In *International Conference on Machine Learning*, pages 2943–2952. PMLR, 2020.
- [3] Chuang Liu, Xueqi Ma, Yinbing Zhan, Liang Ding, Dapeng Tao, Bo Du, Wenbin Hu, and Danilo Mandic. Comprehensive graph gradual pruning for sparse training in graph neural networks. *arXiv preprint arXiv:2207.08629*, 2022.
- [4] Guillaume Bellec, David Kappel, Wolfgang Maass, and Robert Legenstein. Deep rewiring: Training very sparse deep networks. *International Conference on Learning Representations (ICLR)*, 2018.
- [5] Tim Dettmers and Luke Zettlemoyer. Sparse networks from scratch: Faster training without losing performance. *arXiv preprint arXiv:1907.04840*, 2019.
- [6] Shiwei Liu, Lu Yin, Decebal Constantin Mocanu, and Mykola Pechenizkiy. Do we actually need dense over-parameterization? in-time over-parameterization in sparse training. In *International Conference on Machine Learning*, pages 6989–7000. PMLR, 2021.
- [7] Ozan Özdenizci and Robert Legenstein. Training adversarially robust sparse networks via bayesian connectivity sampling. In *International Conference on Machine Learning*, pages 8314–8324. PMLR, 2021.
- [8] Rong Yu and Peichun Li. Toward resource-efficient federated learning in mobile edge computing. *IEEE Network*, 35(1):148–155, 2021.
- [9] Johanna Rock, Wolfgang Roth, Mate Toth, Paul Meissner, and Franz Pernkopf. Resource-efficient deep neural networks for automotive radar interference mitigation. *IEEE Journal of Selected Topics in Signal Processing*, 15(4):927–940, 2021.
- [10] Clayton Frederick Souza Leite and Yu Xiao. Optimal sensor channel selection for resource-efficient deep activity recognition. In *Proceedings of the 20th International Conference on Information Processing in Sensor Networks (co-located with CPS-IoT Week 2021)*, pages 371–383, 2021.
- [11] Torsten Hoefler, Dan Alistarh, Tal Ben-Nun, Nikoli Dryden, and Alexandra Peste. Sparsity in deep learning: Pruning and growth for efficient inference and training in neural networks. *arXiv preprint arXiv:2102.00554*, 2021.
- [12] Laura Graesser, Utku Evci, Erich Elsen, and Pablo Samuel Castro. The state of sparse training in deep reinforcement learning. In *International Conference on Machine Learning*, pages 7766–7792. PMLR, 2022.
- [13] Vikash Sehwal, Shiqi Wang, Prateek Mittal, and Suman Jana. Hydra: Pruning adversarially robust neural networks. *Advances in Neural Information Processing Systems*, 33:19655–19666, 2020.
- [14] Brian Bartoldson, Ari Morcos, Adrian Barbu, and Gordon Erlebacher. The generalization-stability tradeoff in neural network pruning. *Advances in Neural Information Processing Systems*, 33:20852–20864, 2020.
- [15] Xia Xiao, Zigeng Wang, and Sanguthevar Rajasekaran. Autoprune: Automatic network pruning by regularizing auxiliary parameters. *Advances in neural information processing systems*, 32, 2019.

- [16] Shaokai Ye, Kaidi Xu, Sijia Liu, Hao Cheng, Jan-Henrik Lambrechts, Huan Zhang, Aojun Zhou, Kaisheng Ma, Yanzhi Wang, and Xue Lin. Adversarial robustness vs. model compression, or both? In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 111–120, 2019.
- [17] Souvik Kundu, Mahdi Nazemi, Peter A Beerel, and Massoud Pedram. DNR: A tunable robust pruning framework through dynamic network rewiring of DNNs. In *Proceedings of the 26th Asia and South Pacific Design Automation Conference*, pages 344–350, 2021.
- [18] Difan Zou, Pan Xu, and Quanquan Gu. Subsampled stochastic variance-reduced gradient Langevin dynamics. In *International Conference on Uncertainty in Artificial Intelligence*, 2018.
- [19] Changyou Chen, Wenlin Wang, Yizhe Zhang, Qinliang Su, and Lawrence Carin. A convergence analysis for a class of practical variance-reduction stochastic gradient MCMC. *Science China Information Sciences*, 62(1):1–13, 2019.
- [20] Eduard Gorbunov, Filip Hanzely, and Peter Richtárik. A unified theory of SGD: Variance reduction, sampling, quantization and coordinate descent. In *International Conference on Artificial Intelligence and Statistics*, pages 680–690. PMLR, 2020.
- [21] Kumar Avinava Dubey, Sashank J Reddi, Sinead A Williamson, Barnabas Poczos, Alexander J Smola, and Eric P Xing. Variance reduction in stochastic gradient Langevin dynamics. *Advances in neural information processing systems*, 29:1154–1162, 2016.
- [22] Niladri Chatterji, Nicolas Flammarion, Yian Ma, Peter Bartlett, and Michael Jordan. On the theory of variance reduction for stochastic gradient Monte Carlo. In *International Conference on Machine Learning*, pages 764–773. PMLR, 2018.
- [23] Yan Li, Ethan Fang, Huan Xu, and Tuo Zhao. Implicit bias of gradient descent based adversarial training on separable data. In *International Conference on Learning Representations*, 2020.
- [24] Itay Hubara, Brian Chmiel, Moshe Isard, Ron Banner, Joseph Naor, and Daniel Soudry. Accelerated sparse neural training: A provable and efficient method to find n: m transposable masks. *Advances in Neural Information Processing Systems*, 34, 2021.
- [25] Beidi Chen, Tri Dao, Kaizhao Liang, Jiaming Yang, Zhao Song, Atri Rudra, and Christopher Re. Pixelated butterfly: Simple and efficient sparse training for neural network models. *arXiv preprint arXiv:2112.00029*, 2021.
- [26] Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural networks. *International Conference on Learning Representations (ICLR)*, 2019.
- [27] Hesham Mostafa and Xin Wang. Parameter efficient training of deep convolutional neural networks by dynamic sparse reparameterization. In *International Conference on Machine Learning*, pages 4646–4655. PMLR, 2019.
- [28] Siddhant Jayakumar, Razvan Pascanu, Jack Rae, Simon Osindero, and Erich Elsen. Top-kast: Top-k always sparse training. *Advances in Neural Information Processing Systems*, 33:20744–20754, 2020.
- [29] Xiao Zhou, Weizhong Zhang, Hang Xu, and Tong Zhang. Effective sparsification of neural networks with global sparsity constraint. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3599–3608, 2021.
- [30] Jonathan Schwarz, Siddhant Jayakumar, Razvan Pascanu, Peter E Latham, and Yee Teh. Powerpropagation: A sparsity inducing weight reparameterisation. *Advances in Neural Information Processing Systems*, 34:28889–28903, 2021.

- [31] Shaoyi Huang, Bowen Lei, Dongkuan Xu, Hongwu Peng, Yue Sun, Mimi Xie, and Caiwen Ding. Dynamic sparse training via balancing the exploration-exploitation trade-off. *arXiv preprint arXiv:2211.16667*, 2022.
- [32] Bowen Lei, Ruqi Zhang, Dongkuan Xu, and Bani Mallick. Calibrating the rigged lottery: Making all tickets reliable. *arXiv preprint arXiv:2302.09369*, 2023.
- [33] Shaoyi Huang, Haowen Fang, Kaleel Mahmood, Bowen Lei, Nuo Xu, Bin Lei, Yue Sun, Dongkuan Xu, Wujie Wen, and Caiwen Ding. Neurogenesis dynamics-inspired spiking neural network training acceleration. *arXiv preprint arXiv:2304.12214*, 2023.
- [34] Adnan Siraj Rakin, Zhezhi He, Li Yang, Yanzhi Wang, Liqiang Wang, and Deliang Fan. Robust sparse regularization: Simultaneously optimizing neural network robustness and compactness. *arXiv preprint arXiv:1905.13074*, 2019.
- [35] Naveed Akhtar and Ajmal Mian. Threat of adversarial attacks on deep learning in computer vision: A survey. *Ieee Access*, 6:14410–14430, 2018.
- [36] Hongwu Peng, Shaoyi Huang, Tong Geng, Ang Li, Weiwen Jiang, Hang Liu, Shusen Wang, and Caiwen Ding. Accelerating transformer-based deep learning models on fpgas using column balanced block pruning. In *2021 22nd International Symposium on Quality Electronic Design (ISQED)*, pages 142–148. IEEE, 2021.
- [37] Panjie Qi, Edwin Hsing-Mean Sha, Qingfeng Zhuge, Hongwu Peng, Shaoyi Huang, Zhenglun Kong, Yuhong Song, and Bingbing Li. Accelerating framework of transformer by hardware design and model compression co-optimization. In *2021 IEEE/ACM International Conference On Computer Aided Design (ICCAD)*, pages 1–9. IEEE, 2021.
- [38] Nicolas Roux, Mark Schmidt, and Francis Bach. A stochastic gradient method with an exponential convergence rate for finite training sets. *Advances in neural information processing systems*, 25, 2012.
- [39] Rie Johnson and Tong Zhang. Accelerating stochastic gradient descent using predictive variance reduction. *Advances in neural information processing systems*, 26:315–323, 2013.
- [40] Aaron Defazio, Francis Bach, and Simon Lacoste-Julien. SAGA: A fast incremental gradient method with support for non-strongly convex composite objectives. *Advances in neural information processing systems*, 27, 2014.
- [41] Lin Xiao and Tong Zhang. A proximal stochastic gradient method with progressive variance reduction. *SIAM Journal on Optimization*, 24(4):2057–2075, 2014.
- [42] Fanhua Shang, Kaiwen Zhou, Hongying Liu, James Cheng, Ivor W Tsang, Lijun Zhang, Dacheng Tao, and Licheng Jiao. VR-SGD: A simple stochastic variance reduction method for machine learning. *IEEE Transactions on Knowledge and Data Engineering*, 32(1):188–202, 2018.
- [43] Robert M Gower, Mark Schmidt, Francis Bach, and Peter Richtárik. Variance-reduced methods for machine learning. *Proceedings of the IEEE*, 108(11):1968–1983, 2020.
- [44] Lam M Nguyen, Jie Liu, Katya Scheinberg, and Martin Takáč. Sarah: A novel method for machine learning problems using stochastic recursive gradient. In *International Conference on Machine Learning*, pages 2613–2621. PMLR, 2017.
- [45] Cong Fang, Chris Junchi Li, Zhouchen Lin, and Tong Zhang. Spider: Near-optimal non-convex optimization via stochastic path-integrated differential estimator. *Advances in Neural Information Processing Systems*, 31, 2018.
- [46] Ashok Cutkosky and Francesco Orabona. Momentum-based variance reduction in non-convex SGD. *Advances in neural information processing systems*, 32, 2019.

- [47] El Mahdi Chayti and Sai Praneeth Karimireddy. Optimization with access to auxiliary information. *arXiv preprint arXiv:2206.00395*, 2022.
- [48] Xiao Zhou, Weizhong Zhang, Zonghao Chen, Shizhe Diao, and Tong Zhang. Efficient neural network training via forward and backward propagation sparsification. *Advances in Neural Information Processing Systems*, 34:15216–15229, 2021.
- [49] Zeyuan Allen-Zhu and Elad Hazan. Variance reduction for faster non-convex optimization. In *International conference on machine learning*, pages 699–707. PMLR, 2016.
- [50] Jack Baker, Paul Fearnhead, Emily B Fox, and Christopher Nemeth. Control variates for stochastic gradient MCMC. *Statistics and Computing*, 29(3):599–615, 2019.
- [51] Melih Elibol, Lihua Lei, and Michael I Jordan. Variance reduction with sparse gradients. *arXiv preprint arXiv:2001.09623*, 2020.
- [52] Wei Deng, Qi Feng, Georgios Karagiannis, Guang Lin, and Faming Liang. Accelerating convergence of replica exchange stochastic gradient MCMC via variance reduction. *arXiv preprint arXiv:2010.01084*, 2020.
- [53] Ning Qian. On the momentum term in gradient descent learning algorithms. *Neural networks*, 12(1):145–151, 1999.
- [54] Sebastian Ruder. An overview of gradient descent optimization algorithms. *arXiv preprint arXiv:1609.04747*, 2016.
- [55] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [56] Sashank J Reddi, Ahmed Hefny, Suvrit Sra, Barnabás Póczos, and Alex Smola. Stochastic variance reduction for nonconvex optimization. In *International conference on machine learning*, pages 314–323. PMLR, 2016.
- [57] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. *Technical report, University of Toronto*, 2009.
- [58] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y Ng. Reading digits in natural images with unsupervised feature learning. 2011.
- [59] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115(3):211–252, 2015.
- [60] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *International Conference on Learning Representations (ICLR)*, 2015.
- [61] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [62] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. *British Machine Vision Conference*, 2016.
- [63] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. In *International Conference on Learning Representations (ICLR)*, 2018.
- [64] Hongyang Zhang, Yaodong Yu, Jiantao Jiao, Eric Xing, Laurent El Ghaoui, and Michael Jordan. Theoretically principled trade-off between robustness and accuracy. In *International Conference on Machine Learning*, pages 7472–7482. PMLR, 2019.

- [65] Utku Evci, Yani Ioannou, Cem Keskin, and Yann Dauphin. Gradient flow in sparse neural networks and how lottery tickets win. In *Proceedings of the AAAI conference on artificial intelligence*, volume 36, pages 6577–6586, 2022.
- [66] Varun Sundar and Rajat Vadiraj Dwaraknath. [reproducibility report] rigging the lottery: Making all tickets winners. *arXiv preprint arXiv:2103.15767*, 2021.
- [67] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations (ICLR)*, 2019.
- [68] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE international conference on computer vision*, pages 1026–1034, 2015.
- [69] Mahdi Nikdan, Tommaso Pegolotti, Eugenia Iofinova, Eldar Kurtic, and Dan Alistarh. Sparse-prop: efficient sparse backpropagation for faster training of neural networks at the edge. In *International Conference on Machine Learning*, pages 26215–26227. PMLR, 2023.
- [70] Chao Fang, Wei Sun, Aojun Zhou, and Zhongfeng Wang. Efficient n: M sparse dnn training using algorithm, architecture, and dataflow co-design. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, 2023.
- [71] Lei Zhang, Jie Zhang, Bowen Lei, Subhabrata Mukherjee, Xiang Pan, Bo Zhao, Caiwen Ding, Yao Li, and Dongkuan Xu. Accelerating dataset distillation via model augmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11950–11959, 2023.
- [72] Dongyao Zhu, Bowen Lei, Jie Zhang, Yanbo Fang, Yiqun Xie, Ruqi Zhang, and Dongkuan Xu. Rethinking data distillation: Do not overlook calibration. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4935–4945, 2023.
- [73] Cheng Wang, Jiacheng Sun, Zhenhua Dong, Ruixuan Li, and Rui Zhang. Gradient matching for categorical data distillation in ctr prediction. In *Proceedings of the 17th ACM Conference on Recommender Systems*, pages 161–170, 2023.
- [74] Jiangyuan Li, Thanh Nguyen, Chinmay Hegde, and Ka Wai Wong. Implicit sparse regularization: The impact of depth and early stopping. *Advances in Neural Information Processing Systems*, 34:28298–28309, 2021.
- [75] Hung-Hsu Chou, Johannes Maly, and Holger Rauhut. More is less: inducing sparsity via overparameterization. *Information and Inference: A Journal of the IMA*, 12(3):iaad012, 2023.
- [76] Jiangyuan Li, Thanh V Nguyen, Chinmay Hegde, and Raymond KW Wong. Implicit regularization for group sparsity. *arXiv preprint arXiv:2301.12540*, 2023.
- [77] Zhiyao Li, Jiaxiang Li, Taijie Chen, Dimin Niu, Hongzhong Zheng, Yuan Xie, and Mingyu Gao. Spada: Accelerating sparse matrix multiplication with adaptive dataflow. In *Proceedings of the 28th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 2*, pages 747–761, 2023.

A. Appendix: Theoretical Proof of Convergence Rate

In this section, we provide detailed proof of the convergence rate of our AGENT method. We start with some assumptions on which we will give some useful lemmas. Then, we will establish the convergence rate of our AGENT method based on these lemmas.

A.1. Algorithm Reformulation

We reformulate our Adaptive Gradient Correction (AGENT) into a math-friendly version that is shown in Algorithm 2.

Algorithm 2 Adaptive Gradient Correction

Input: Initialize θ_0^0 and $c_{-1} = 0$, set the number of epochs S , epoch length m , step sizes η_t , scaling parameter γ , and smoothing factor α

for $s = 0$ **to** $S - 1$ **do**

$\tilde{\theta} = \theta_0^s$

$\tilde{g} = (\sum_{i=1}^N \nabla l(\mathbf{x}_i; \tilde{\theta})) / N$

Calculate \tilde{c}_s^* via Eq. (4)

$\tilde{c}_s = (1 - \alpha)\tilde{c}_{s-1} + \alpha\tilde{c}_s^*$

$c_s = \gamma\tilde{c}_s$

for $t = 0$ **to** $m - 1$ **do**

Sample a mini-batch data \mathbf{B}_t with size n

$\theta_{t+1}^s = \theta_t^s - \eta_t \left(\frac{1}{n} \sum_{i \in \mathbf{B}_t} (g_i(\theta_t^s) - c_s \cdot g_i(\tilde{\theta})) + c_s \cdot \tilde{g} \right)$

end for

$\theta_0^{s+1} = \theta_m^s$

end for

Output: Iterates θ_π chosen uniformly random from $\{\{\theta_t^s\}_{t=0}^{m-1}\}_{s=0}^{S-1}$

A.2. Assumptions

L-smooth: A differentiable function $l : \mathbb{R}^n \rightarrow \mathbb{R}$ is said to be L-smooth if for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ is satisfies $\|\nabla l(\mathbf{x}) - \nabla l(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|$. An equivalent definition is for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$:

$$-\frac{L}{2}\|\mathbf{x} - \mathbf{y}\|^2 \leq l(\mathbf{x}) - l(\mathbf{y}) - \langle \nabla l(\mathbf{x}), \mathbf{x} - \mathbf{y} \rangle \leq \frac{L}{2}\|\mathbf{x} - \mathbf{y}\|^2$$

σ -bounded: We say function l has a σ -bounded gradient if $\|\nabla l_i(\mathbf{x})\| \leq \sigma$ for all $i \in [N]$ and $\mathbf{x} \in \mathbb{R}^n$

A.3. Analysis framework

Under the above assumptions, we are ready to analyze the convergence rate of AGENT in **Algorithm 2**. To introduce the convergence analysis more clearly, we provide a brief analytical framework for our proof.

- First, we need to show that the variance of our gradient estimator is smaller than that of minibatch SVRG under proper choice of c_s . Since the gradient estimator of both AGENT and minibatch SVRG are unbiased estimators in standard training, we only need to show that our bound $E[\|\mathbf{u}_t\|^2]$ is smaller than minibatch SVRG. (See in Lemma 1)
- Based on above fact, we next apply the Lyapunov function to prove the convergence rate of AGENT in one arbitrary epoch. (See in Lemma 3)
- Then, we extend our previous results to the entire epoch (from 0 to S -th epoch) and derive the convergence rate of the output θ_π of **Algorithm 2**. (See in Lemma 4)

- Finally, we compare the convergence rate of our AGENT with that of minibatch SVRG. Setting the parameters in Lemma 4 according to the actual situation of sparse learning, we obtain a bound that is more stringent than minibatch SVRG.

A.4. Lemma

We first denote step length $\eta_t = N \cdot h_t$. Since we mainly focus on a single epoch, we drop the superscript s and denote $\mathbf{u}_t = \frac{1}{n} \sum_{i \in \mathbf{B}_t} \left(g_i(\boldsymbol{\theta}_t) - c \cdot g_i(\tilde{\boldsymbol{\theta}}) \right) + c \cdot \tilde{g}$ which is the gradient estimator in our algorithm and $\boldsymbol{\tau}_t = \frac{1}{n} \sum_{i \in \mathbf{B}_t} \left(g_i(\boldsymbol{\theta}_t) - c \cdot g_i(\tilde{\boldsymbol{\theta}}) \right)$, then lines the update procedure in **Algorithm 2** can be replaced with $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta_t \cdot \mathbf{u}_t$

A.4.1. Lemma 1

For the \mathbf{u}_t defined above and function l is a L -smooth, λ -strongly convex function with σ -bounded gradient, then we have the following results:

$$\mathbb{E} [\|\mathbf{u}_t\|^2] \leq 2\mathbb{E} [\|g(\boldsymbol{\theta}_t)\|^2] + \frac{4c^2L^2}{n} \mathbb{E} [\|\boldsymbol{\theta}_t - \tilde{\boldsymbol{\theta}}\|^2] + \frac{4(1-c)^2}{n} \sigma^2 \quad (6)$$

Proof:

$$\begin{aligned} \mathbb{E} [\|\mathbf{u}_t\|^2] &= \mathbb{E} [\|\boldsymbol{\tau}_t + c \cdot \tilde{g}\|^2] = \mathbb{E} [\|\boldsymbol{\tau}_t + c \cdot \tilde{g} - g(\boldsymbol{\theta}_t) + g(\boldsymbol{\theta}_t)\|^2] \\ &\leq 2\mathbb{E} [\|g(\boldsymbol{\theta}_t)\|^2] + 2\mathbb{E} [\|\boldsymbol{\tau}_t - \mathbb{E}[\boldsymbol{\tau}_t]\|^2] \leq 2\mathbb{E} [\|g(\boldsymbol{\theta}_t)\|^2] + \frac{2}{n} \mathbb{E} [\boldsymbol{\tau}_t^2] \\ &= 2\mathbb{E} [\|g(\boldsymbol{\theta}_t)\|^2] + \frac{2}{n} \mathbb{E} [\|c(g_i(\boldsymbol{\theta}_t) - g_i(\tilde{\boldsymbol{\theta}})) + (1-c)g_i(\boldsymbol{\theta}_t)\|^2] \\ &\leq 2\mathbb{E} [\|g(\boldsymbol{\theta}_t)\|^2] + \frac{4}{n} \mathbb{E} [\|c(g_i(\boldsymbol{\theta}_t) - g_i(\tilde{\boldsymbol{\theta}}))\|^2] + \frac{4(1-c)^2}{n} \mathbb{E} [\|g_i(\boldsymbol{\theta}_t)\|^2] \\ &\leq 2\mathbb{E} [\|g(\boldsymbol{\theta}_t)\|^2] + \frac{4c^2L^2}{n} \mathbb{E} [\|\boldsymbol{\theta}_t - \tilde{\boldsymbol{\theta}}\|^2] + \frac{4(1-c)^2}{n} \sigma^2 \end{aligned}$$

The first and third inequality are because $\|a+b\|^2 \leq 2\|a\|^2 + 2\|b\|^2$, the second inequality follows the $\mathbb{E} [\|\boldsymbol{\tau} - \mathbb{E}[\boldsymbol{\tau}]\|^2] \leq \mathbb{E} [\|\boldsymbol{\tau}\|^2]$ and the last inequality follows the L -smoothness and σ -bounded of function l_i .

Remark 5. Compared with the gradient estimator of minibatch SVRG, the bound of $\mathbb{E}[\|\mathbf{u}_t\|^2]$ is smaller when L is large, σ is relatively small and c is properly chosen.

A.4.2. Lemma 2

$$\mathbb{E} [l(\boldsymbol{\theta}_{t+1})] \leq \mathbb{E} \left[l(\boldsymbol{\theta}_t) + \eta_t \|g(\boldsymbol{\theta}_t)\|^2 + \frac{L\eta_t^2}{2} \|\mathbf{u}_t\|^2 \right] \quad (7)$$

Proof:

By the L -smoothness of function l , we have

$$\mathbb{E} [l(\boldsymbol{\theta}_{t+1})] \leq \mathbb{E} \left[l(\boldsymbol{\theta}_t) + \langle g(\boldsymbol{\theta}_t), \boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_t \rangle + \frac{L}{2} \|\boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_t\|^2 \right]$$

By the update procedure in algorithm 2 and unbiasedness, the right hand side can further upper bounded by

$$\mathbb{E} \left[l(\boldsymbol{\theta}_t) + \eta_t \|g(\boldsymbol{\theta}_t)\|^2 + \frac{L\eta_t^2}{2} \|\mathbf{u}_t\|^2 \right]$$

A.4.3. Lemma 3

For $b_t, b_{t+1}, \zeta_t > 0$ and b_t and b_{t+1} have the following relationship

$$b_t = b_{t+1} \left(1 + \eta_t \zeta_t + \frac{4c^2 \eta_t^2 L^2}{n} \right) + 2 \frac{c^2 \eta_t^2 L^3}{n}$$

and define

$$\Phi_t := \eta_t - \frac{b_{t+1} \eta_t}{\zeta_t} - \eta_t^2 L - 2b_{t+1} \eta_t^2$$

$$\Psi_t := \mathbb{E} \left[l(\boldsymbol{\theta}_t) + b_t \|\boldsymbol{\theta}_t - \tilde{\boldsymbol{\theta}}\|^2 \right] \quad (8)$$

η_t, ζ_t and b_{t+1} can be chosen such that $\Phi_t > 0$. Then the x_t in Algorithm 1 have the bound:

$$\mathbb{E}[\|g(\boldsymbol{\theta}_t)\|^2] \leq \frac{\Psi_t - \Psi_{t+1} + \frac{2(L\eta_t^2 + 2b_{t+1}\eta_t^2)(1-c)^2}{n} \sigma^2}{\Phi_t}$$

Proof:

We apply the Lyapunov function

$$\Psi_t = \mathbb{E} \left[l(\boldsymbol{\theta}_t) + b_t \|\boldsymbol{\theta}_t - \tilde{\boldsymbol{\theta}}\|^2 \right]$$

Then we need to bound $\|\boldsymbol{\theta}_t - \tilde{\boldsymbol{\theta}}\|$

$$\begin{aligned} \mathbb{E} \left[\|\boldsymbol{\theta}_{t+1} - \tilde{\boldsymbol{\theta}}\|^2 \right] &= \mathbb{E} \left[\|\boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_t + \boldsymbol{\theta}_t - \tilde{\boldsymbol{\theta}}\|^2 \right] \\ &= \mathbb{E} \left[\|\boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_t\|^2 + \|\boldsymbol{\theta}_t - \tilde{\boldsymbol{\theta}}\|^2 + 2\langle \boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_t, \boldsymbol{\theta}_t - \tilde{\boldsymbol{\theta}} \rangle \right] \\ &= \mathbb{E} \left[\eta_t^2 \|\mathbf{u}_t\|^2 + \|\boldsymbol{\theta}_t - \tilde{\boldsymbol{\theta}}\|^2 \right] - 2\eta_t \mathbb{E} \left[\langle g(\boldsymbol{\theta}_t), \boldsymbol{\theta}_t - \tilde{\boldsymbol{\theta}} \rangle \right] \\ &\leq \mathbb{E}[\eta_t^2 \|\mathbf{u}_t^{s+1}\|^2 + \|\boldsymbol{\theta}_t - \tilde{\boldsymbol{\theta}}\|^2] + 2\eta_t \mathbb{E} \left[\frac{1}{2\zeta_t} \|g(\boldsymbol{\theta}_t)\| + \frac{\zeta_t}{2} \|\boldsymbol{\theta}_t - \tilde{\boldsymbol{\theta}}\|^2 \right] \end{aligned} \quad (9)$$

The third equality is due to the unbiasedness of the update and the last inequality follows Cauchy-Schwarz and Young's inequality. Plugging Equation (6), Equation (7), and Equation (9) into Equation (8), we can get the following bound:

$$\begin{aligned} \Psi_{t+1} &\leq \mathbb{E} [l(\boldsymbol{\theta}_t)] + \left(b_{t+1} \left(1 + \eta_t \zeta_t + \frac{4c^2 \eta_t^2 L^2}{n} \right) + \frac{2c^2 \eta_t^2 L^3}{n} \right) \mathbb{E} [\|\boldsymbol{\theta}_t - \tilde{\boldsymbol{\theta}}\|^2] \\ &\quad - \left(\eta_t - \frac{b_{t+1} \eta_t}{\zeta_t} - L\eta_t^2 - 2b_{t+1} \eta_t^2 \right) \mathbb{E} [\|g(\boldsymbol{\theta}_t)\|^2] + 4 \left(\frac{L\eta_t^2}{2} + b_{t+1} \eta_t^2 \right) \frac{(1-c)^2}{n} \sigma^2 \\ &= \Psi_t - \left(\eta_t - \frac{b_{t+1} \eta_t}{\zeta_t} - L\eta_t^2 - 2b_{t+1} \eta_t^2 \right) \mathbb{E} [\|g(\boldsymbol{\theta}_t)\|^2] + 4 \left(\frac{L\eta_t^2}{2} + b_{t+1} \eta_t^2 \right) \frac{(1-c)^2}{n} \sigma^2 \end{aligned}$$

A.4.4. Lemma 4

Now we consider the effect of epoch and use s to denote the epoch number. Let $b_m^s = 0, \eta_t^s = \eta, \zeta_t^s = \zeta$ and $b_t^s = b_{t+1}^s \left(1 + \eta \zeta + \frac{4c^2 \eta L^2}{n} \right) + 2 \frac{c^2 \eta^2 L^2}{n}, \Phi_t^s = \eta - \frac{b_{t+1}^s \eta}{\zeta_t} - \eta^2 L - 2b_{t+1}^s \eta^2$. Define $\phi := \min_{t,s} \Phi_t^s$. Then we can conclude that:

$$\mathbb{E}[\|g(\boldsymbol{\theta}_\pi)\|^2] \leq \frac{l(\theta_0) - l(\theta_*)}{T\phi} + \sum_{s=0}^{S-1} \sum_{t=0}^{m-1} \frac{2(L + 2b_{t+1}^s)(1 - c_s)^2 \eta^2 \sigma^2}{Tn\phi}$$

Proof:

Under the condition of $\eta_t^s = \eta$, we apply telescoping sum on **Lemma 3**, then we will get:

$$\sum_{t=1}^{m-1} \mathbb{E}[\|g(\boldsymbol{\theta}_t^s)\|^2] \leq \frac{\Psi_0^s - \Psi_m^s}{\phi} + \sum_{t=0}^{m-1} \frac{2(L + 2b_{t+1}^s)(1 - c_s)^2 \eta^2 \sigma^2}{n\phi}$$

From previous definition, we know $\Psi_0^s = l(\tilde{\boldsymbol{\theta}}^s)$, $\Psi_m^s = l(\tilde{\boldsymbol{\theta}}^{s+1})$ and plugging into previous equation, we obtain:

$$\sum_{t=1}^{m-1} \mathbb{E}[\|g(\boldsymbol{\theta}_t^s)\|^2] \leq \frac{l(\tilde{\boldsymbol{\theta}}^s) - l(\tilde{\boldsymbol{\theta}}^{s+1})}{\phi} + \sum_{t=0}^{m-1} \frac{2(L + 2b_{t+1}^s)(1 - c_s)^2 \eta^2 \sigma^2}{n\phi}$$

Take summation over all the epochs and using the fact that $\tilde{\boldsymbol{\theta}}^0 = \boldsymbol{\theta}_0$, $l(\tilde{\boldsymbol{\theta}}^S) \leq l(\boldsymbol{\theta}_*)$ we immediately obtain:

$$\frac{1}{T} \sum_{s=0}^{S-1} \sum_{t=1}^{m-1} \mathbb{E}[\|g(\boldsymbol{\theta}_t^s)\|^2] \leq \frac{l(\boldsymbol{\theta}_0) - l(\boldsymbol{\theta}_*)}{\phi} + \sum_{s=0}^{S-1} \sum_{t=0}^{m-1} \frac{2(L + 2b_{t+1}^s)(1 - c_s)^2 \eta^2 \sigma^2}{Tn\phi} \quad (10)$$

A.5. Theorem

A.5.1. Theorem 1

Define $\xi_s = \sum_{t=0}^{m-1} (L + 2b_{t+1}^s)$ and $\xi := \min_s \xi_s$. Let $\eta = \frac{\mu n}{LN^\alpha}$ ($0 < \mu < 1$) and ($0 < \alpha \leq 1$), $\zeta = \frac{L}{N^{\alpha/2}}$ and $m = \frac{N^{\frac{3\alpha}{2}}}{\mu n}$. Then there exists constant $\nu, \mu, \alpha, \kappa > 0$ such that $\phi \geq \frac{n\nu}{LN^\alpha}$ and $\xi \leq \kappa L$. Then $\mathbb{E}[\|g(\boldsymbol{\theta}_\pi)\|^2]$ can be future bounded by:

$$\mathbb{E}[\|g(\boldsymbol{\theta}_\pi)\|^2] \leq \frac{(l(\boldsymbol{\theta}_0) - l(\boldsymbol{\theta}_*))LN^\alpha}{Tn\nu} + \frac{2\kappa\mu^2\sigma^2}{N^\alpha\nu m}$$

Proof:

By applying summation formula of geometric progression on the relation $b_t^s = b_{t+1}^s(1 + \eta_t\zeta_t + \frac{4c_s^2\eta_t^2L^2}{n}) + 2\frac{c_s^2\eta_t^2L^3}{n}$, we have $b_t^s = \frac{2c_s^2\eta^2L^3}{n} \frac{(1+\omega_s)^{m-t}-1}{\omega_s}$ where:

$$\omega_s = \eta\zeta + \frac{4c_s^2\eta^2L}{n} = \frac{\mu n}{N^{\frac{3\alpha}{2}}} + \frac{4c_s^2\mu^2n}{N^{2\alpha}} \leq \frac{(4c_s^2 + 1)\mu n}{N^{\frac{3\alpha}{2}}}$$

This bound holds because $\mu \leq 1$ and $N \geq 1$ and thus $\frac{4c_s^2\mu^2n}{N^{2\alpha}} = \frac{4c_s^2\mu n}{N^{\frac{3\alpha}{2}}} \times \frac{\mu}{N^{\frac{\alpha}{2}}} \leq \frac{4c_s^2\mu n}{N^{\frac{3\alpha}{2}}}$. Using this bound, we obtain:

$$\begin{aligned}
b_0^s &= \frac{2\eta^2 c_s^2 L^3 (1 + \omega_s)^m - 1}{n \omega_s} = \frac{2\mu^2 n c_s^2 L (1 + \omega_s)^m - 1}{N^{2\alpha} \omega_s} \\
&\leq \frac{2\mu n c_s^2 L ((1 + \omega_s)^m - 1)}{N^{\frac{\alpha}{2}} (4c_s + 1)} \leq \frac{2\mu n c_s^2 L \left(\left(1 + \frac{(4c_s^2 + 1)\mu n}{N^{\frac{3\alpha}{2}}}\right)^{\frac{N^{\frac{3\alpha}{2}}}{\mu n}} - 1 \right)}{N^{\frac{\alpha}{2}} (4c_s^2 + 1)} \\
&\leq \frac{2\mu n c_s^2 L (e^{\frac{1}{4c_s^2 + 1}} - 1)}{N^{\frac{\alpha}{2}} (4c_s^2 + 1)}
\end{aligned}$$

The last inequality holds because $(1 + \frac{1}{x})^x$ is a monotone increasing function of x when $x > 0$. Thus $(1 + \frac{(4c_s^2 + 1)\mu n}{N^{\frac{3\alpha}{2}}})^{\frac{N^{\frac{3\alpha}{2}}}{\mu n}} \leq e^{\frac{1}{4c_s^2 + 1}}$ in the third inequality. And we can obtain the lower bound for ϕ

$$\phi = \min_{t,s} \Phi_t^s \geq \min_s (\eta - \frac{b_0^s \eta}{\zeta} - \eta^2 L - 2b_0^s \eta^2) \geq \frac{\nu}{LN^\alpha}$$

The first inequality holds since b_s^t is a decrease function of t . Meanwhile, the second inequality holds because there exists uniform constant ν such that $\nu \geq \mu(1 - \frac{b_0^s \eta}{\zeta} - L\eta - b_0^s \eta)$.

Remark 6. In practice, $b_0^s \approx 0$ because both γ and c_s is both smaller than 0.1 which leads to $\mu(1 - \frac{b_0^s}{\zeta} - L\eta - b_0^s \eta) \approx \mu(1 - L\eta)$ and this value is usually much bigger than the ν^* in the bound of minibatch SVRG.

We need to find the upper bound for ξ

$$\begin{aligned}
\xi_s &= \sum_{t=0}^{m-1} (L + 2b_{t+1}^s) = mL + 2 \sum_{t=0}^{m-1} b_{t+1}^s \\
&= mL + 2 \sum_{t=0}^{m-1} \frac{2c_s^2 \eta^2 L^3 (1 + \omega_s)^{m-t} - 1}{n \omega_s} \\
&= mL + \frac{2c_s^2 \eta^2 L^3}{n \omega_s} \left[\frac{(1 + \omega_s)^{m+1} - (1 + \omega_s)}{\omega_s} - m \right] \\
&\leq mL + \frac{2c_s^2 \eta^2 L^3}{n} \left[\frac{1 + \omega_s}{\omega_s^2} (e^{\frac{1}{4c_s^2 + 1}} - 1) - m \right] \\
&\leq mL + \frac{2c_s^2 L N^\alpha}{n} \left(1 + \frac{\mu n}{N^{3\alpha/2}}\right) (e^{\frac{1}{4c_s^2 + 1}} - 1) - \frac{2c_s^2 \mu^2 n m L}{N^{2\alpha}} \\
&= L \left[\left(1 - \frac{2c_s^2 \mu^2 n L}{N^{2\alpha}}\right) m + \frac{2c_s^2 N^\alpha}{n} \left(1 + \frac{\mu n}{N^{3\alpha/2}}\right) (e^{\frac{1}{4c_s^2 + 1}} - 1) \right]
\end{aligned}$$

The reason why the first inequality holds is explained before and the second inequality holds because $\frac{1+x}{x^2}$ is a monotone decreasing function of x when $x > 0$, $\omega_s = \frac{\mu n}{N^{\frac{3\alpha}{2}}} + \frac{4c_s^2 \mu^2 n}{N^{2\alpha}} \leq \frac{\mu n}{N^{\frac{3\alpha}{2}}}$ and $\eta = \frac{\mu n}{LN^\alpha}$. Then $\xi = \max_s \xi_s \leq \kappa L$ where $\kappa \geq \max_s \left(\left(1 - \frac{2c_s^2 \mu^2 n L}{N^{2\alpha}}\right) m + \frac{2c_s^2 N^\alpha}{n} \left(1 + \frac{\mu n}{N^{3\alpha/2}}\right) (e^{\frac{1}{4c_s^2 + 1}} - 1) \right)$. When $c_s \approx 0$, $\left(1 - \frac{2c_s^2 \mu^2 n L}{N^{2\alpha}}\right) m + \frac{2c_s^2 N^\alpha}{n} \left(1 + \frac{\mu n}{N^{3\alpha/2}}\right) (e^{\frac{1}{4c_s^2 + 1}} - 1) \approx m$.

Now we obtain the lower bound for ϕ and upper bound for ξ , plugging them into equation (10), we will have:

$$\begin{aligned}
\mathbb{E}[\|g(\boldsymbol{\theta}_\pi)\|^2] &\leq \frac{l(\boldsymbol{\theta}_0) - l(\boldsymbol{\theta}_*)}{\phi} + \sum_{s=0}^{S-1} \sum_{t=0}^{m-1} \frac{2(L + 2b_{t+1}^s)(1 - c_s)^2 \eta^2 \sigma^2}{Tn\phi} \\
&\leq \frac{(l(\boldsymbol{\theta}_0) - l(\boldsymbol{\theta}_*))LN^\alpha}{Tn\nu} + \sum_{s=0}^{S-1} \sum_{t=0}^{m-1} \frac{2(L + 2b_{t+1}^s)\eta^2 \sigma^2}{Tn\phi} \\
&\leq \frac{(l(\boldsymbol{\theta}_0) - l(\boldsymbol{\theta}_*))LN^\alpha}{Tn\nu} + \sum_{s=0}^{S-1} \left(\frac{2\eta^2 \sigma^2}{Tn\phi}\right) \sum_{t=0}^{m-1} (L + 2b_{t+1}^s) \\
&\leq \frac{(l(\boldsymbol{\theta}_0) - l(\boldsymbol{\theta}_*))LN^\alpha}{Tn\nu} + \frac{2\kappa\mu^2\sigma^2}{N^\alpha\nu m}
\end{aligned}$$

Remark 7. In our theoretical analysis above, we consider c as a constant in each epoch, which is still consistent with our practical algorithm for the following reasons.

- (i) In our Algorithm 1, \widehat{c}_t^* is actually a fixed constant within each epoch, which can be different in different epochs. Since it is too expensive to compute the exact \widehat{c}_t^* in each iteration, we compute it at the beginning of each epoch and use it as an approximation in the following epoch.
- (ii) As for our proof, we first show the convergence rate of one arbitrary training epoch. In this step, treating c as a constant is aligned with our practical algorithm.
- (iii) Then, when we extend the results of one epoch to the whole epoch, we establish an upper bound for different c in each epoch. Thus, the bound can be applied when c differs across epochs, which enables our theoretical analysis consistent with our practical algorithm.

A.6. Real Case Analysis for Sparse Training

A.6.1. CIFAR-10/100 dataset

In our experiments, we apply both SVRG and AGENT on CIFAR-10 and CIFAR-100 datasets with $\eta = 0.1$, $\gamma = 0.1$, batch size $m = 128$, and in total 50000 training samples. Under this parameter setting, ν and ν^* in **Theorem 1** and **Remark 4** are about 0.1 and 0.06, respectively. While $\frac{2\kappa\mu^2\sigma^2}{N^\alpha\nu m}$ is around 10^{-5} which is negligible so we know AGENT should have a tighter bound than SVRG in this situation which matches with the experimental results show in *Figure 6*.

A.6.2. svhn dataset

Meanwhile, in the SVHN dataset, we train our model with parameters: $\eta = 0.1$, $\gamma = 0.1$, batch size $m = 573$, and sample size $N = 73257$. ν , ν^* equal 0.4 and 0.06 respectively and $\frac{2\kappa\mu^2\sigma^2}{N^\alpha\nu m}$ is around 10^{-4} . Although the second term in **Theorem 1** is bigger. Since ν here is a lot bigger than ν^* which leads to the first term in **Theorem 1** much smaller than that of **Remark 4**. So we still obtain a more stringent bound compared with SVRG which also meets with the outcome presented in *Figure 9*.

B. Additional Experimental Results

We summarize additional experimental results for the BSR-Net-based [7], RigL-based [2], and ITOP-based [6] models.

B.1. Accuracy Comparisons in Different Epochs

Aligned with the main manuscript, we compare the accuracy for a given number of epochs to compare both the speed of convergence and training stability. We first show BSR-Net-based results in this section. Since our approach has faster convergence and does not require a long warm-up period,

the dividing points for the decay scheduler are set to the 50th and 100th epochs. In the manuscript, we also use this schedule for BSR-Net for an accurate comparison. In the Appendix, we include the results using its original schedule. BSR-Net and BSR-Net (ori) represent the results learned using our learning rate schedule and original schedule in [7], respectively. As shown in Figures 6, 7, 8, 9, 10, 11, 12, the blue curves (A-BSR-Net) are always higher than the yellow curves and also much smoother than yellow curves (BSR-Net and BSR-Net (ori)), indicating faster and more stable training when using our proposed A-BSR-Net.

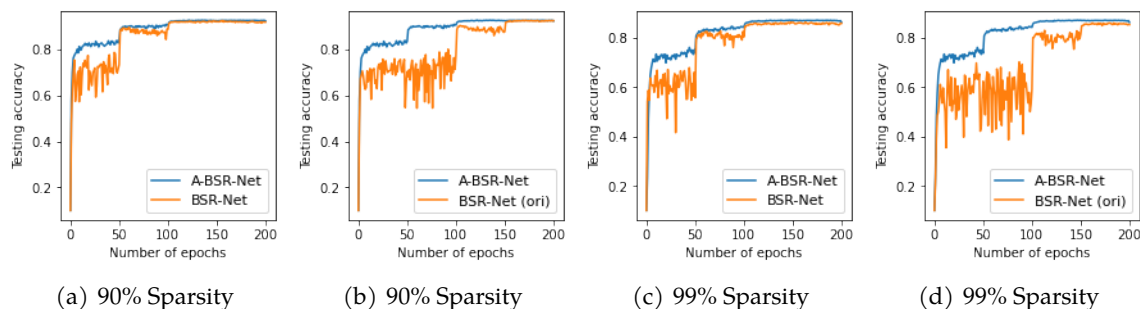


Figure 6: Comparisons (accuracy given the number of epochs) with BSR-Net [7]. We evaluate sparse networks (99% or 90%) learned with natural training on CIFAR-10 using VGG-16.

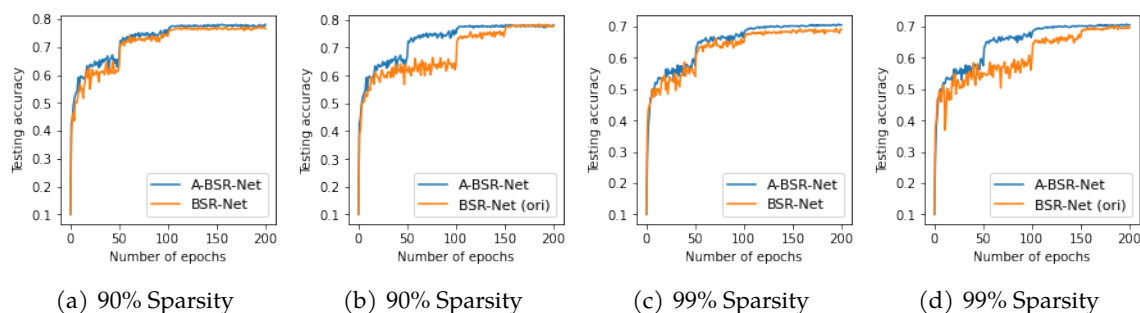


Figure 7: Comparisons (accuracy given the number of epochs) with BSR-Net [7]. We evaluate sparse networks (99% or 90%) learned with adversarial training (objective: AT) on CIFAR-10 using VGG-16.

We also show ITOP-based results on ImageNet-2012. As shown in Figure 13, the red and blue curves represent AGENT + RigL-ITOP and RigL-ITOP on 80% and 90% sparse ResNet-50, respectively. For 80% sparsity, the red curve is above the blue curve, demonstrating the acceleration effect of our AGENT, especially in the early stages. For 90% sparsity, we can see that the red curve is more stable than the blue curve, which shows the stable effect of our AGENT on large data sets. If we use SVRG in this case, we will not only fail to train stably but also slow down the training speed. In contrast, our AGENT can solve the limitation of SVRG.

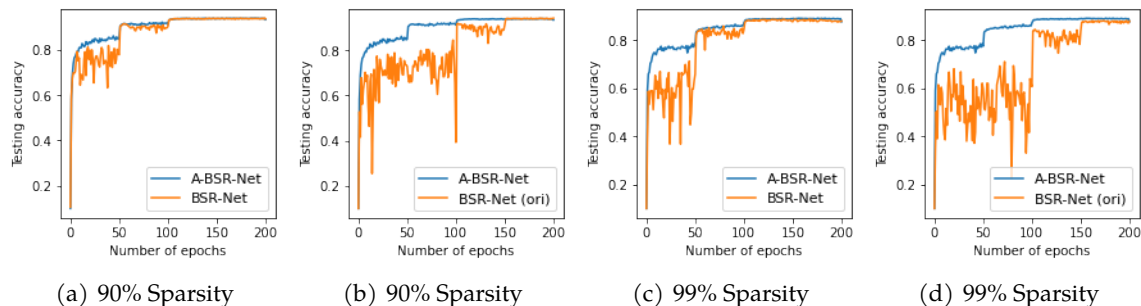


Figure 8: Comparisons (accuracy given the number of epochs) with BSR-Net [7]. We evaluate sparse networks (99% or 90%) learned with natural training on CIFAR-10 using Wide-ResNet-28-4.

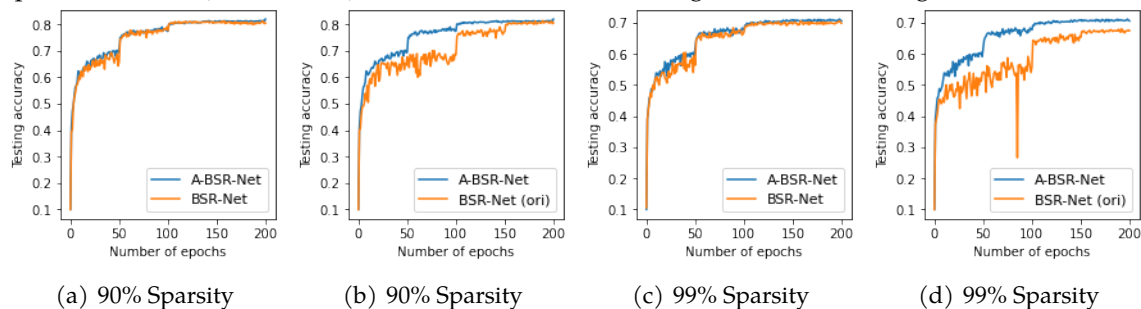


Figure 9: Comparisons (accuracy given the number of epochs) with BSR-Net [7]. We evaluate sparse networks (99% or 90%) learned with adversarial training (objective: AT) on CIFAR-10 using Wide-ResNet-28-4.

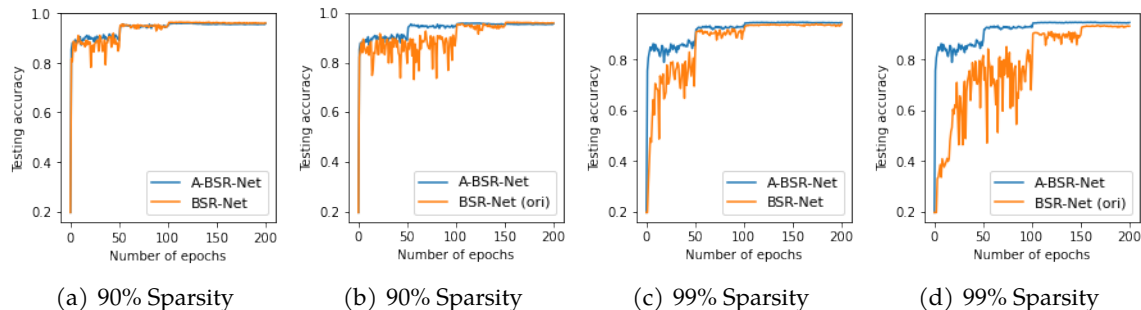


Figure 10: Comparisons (accuracy given the number of epochs) with BSR-Net [7]. We evaluate sparse networks (99% or 90%) learned with natural training on SVHN using VGG-16.

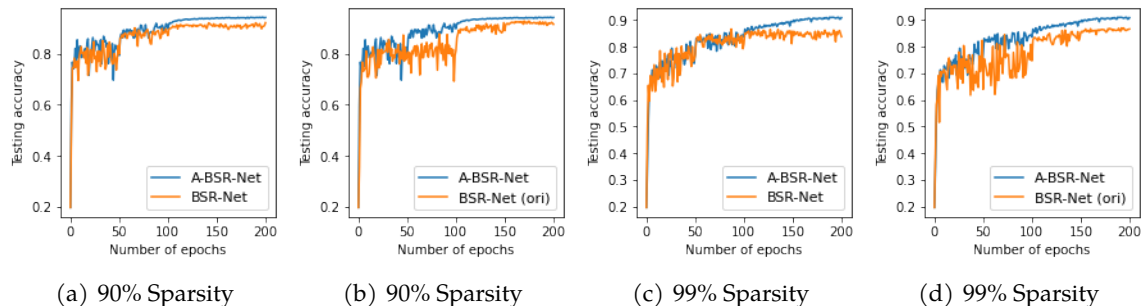
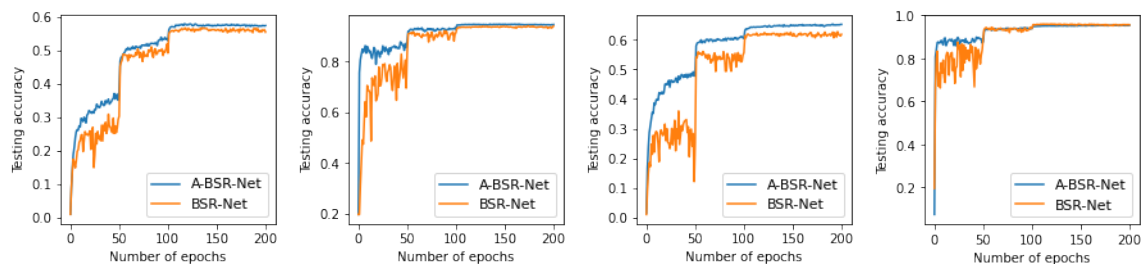
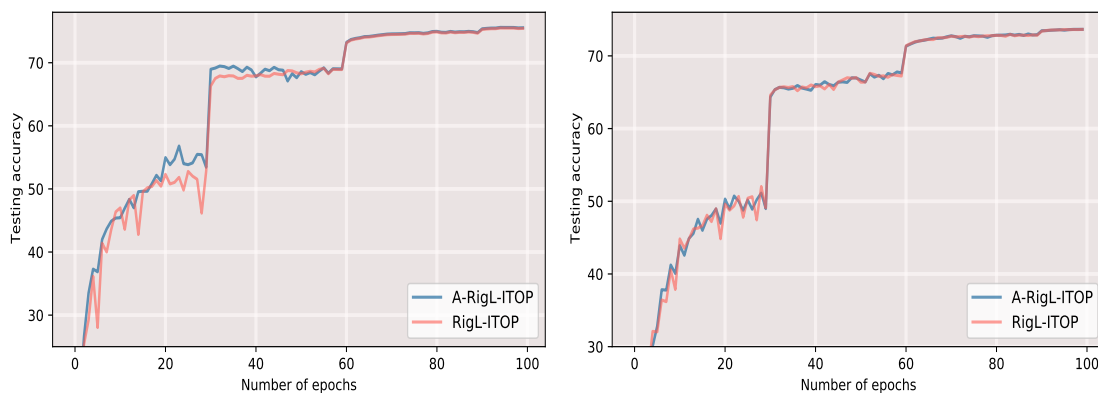


Figure 11: Comparisons (accuracy given the number of epochs) with BSR-Net [7]. We evaluate sparse networks (99% or 90%) learned with adversarial training (objective: TRADES) on SVHN using VGG-16.



(a) CIFAR-100,VGG-16 (b) SVHN,VGG-16 (c) CIFAR-100,WRN-28-4 (d) SVHN,WRN-28-4

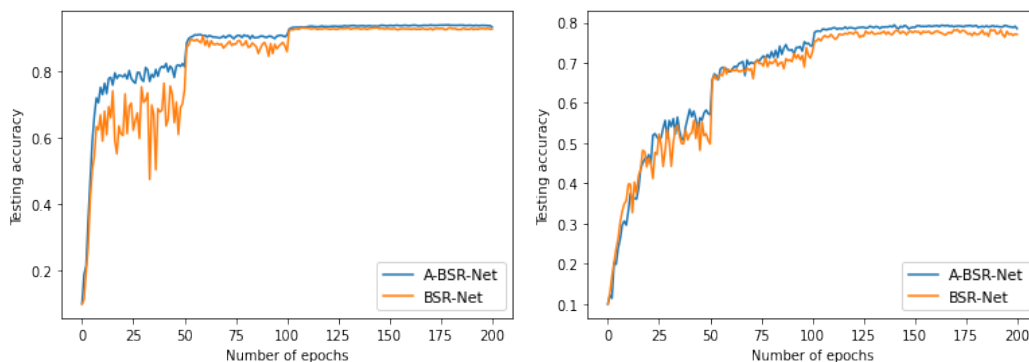
Figure 12: Training curve (accuracy given a number of epochs) of BSR-Net-based models [7]. Sparse networks (99%) are learned in standard setups on (a) CIFAR-100 using VGG-16, (b) SVHN using VGG-16, (c) CIFAR-100 using WRN-28-4, (d) SVHN using WRN-28-4.



(a) 80% Sparsity

(b) 90% Sparsity

Figure 13: Testing accuracy for ITOP-based models at 80% and 90% sparsity on ImageNet-2012. A-RigL-ITOP (blue curves) converges faster than RigL-ITOP (pink curves).



(a) Standard

(b) Adversarial (AT)

Figure 14: Training curve (required epochs to reach given accuracy) of BSR-Net-based models [7]. Dense networks are learned in standard and adversarial setups on CIFAR-10 using VGG-16.

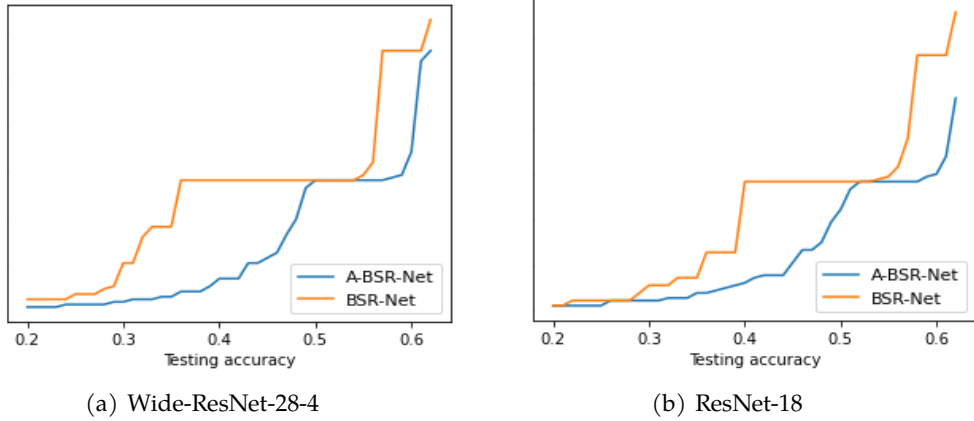


Figure 15: Comparisons (required hours to reach given accuracy. We evaluate sparse networks (99%) learned with natural training on CIFAR-100 using (a) Wide-ResNet-28-4, and (b) ResNet-18.

In Figure 14, we also compare the convergence speed without sparsity. We show a BSR-Net-based result, where the dense network is learned by adversarial training (AT) and standard training on CIFAR-10 using VGG-16. The blue curve of our A-BSR-Net tends to be above the yellow curve of BSR-Net, indicating successful acceleration. This demonstrates the broad applicability of our method.

B.2. Number of Training Epoch Comparisons

We also compare the number of training epochs required to reach the same accuracy in BSR-Net-based results. In Figures 15, 16, 17, 18, 19, 20, 21, the blue curves (A-BSR-Net) are always lower than yellow curves (BSR-Net and BSR-Net (ori)), indicating faster convergence of A-BSR-Net.

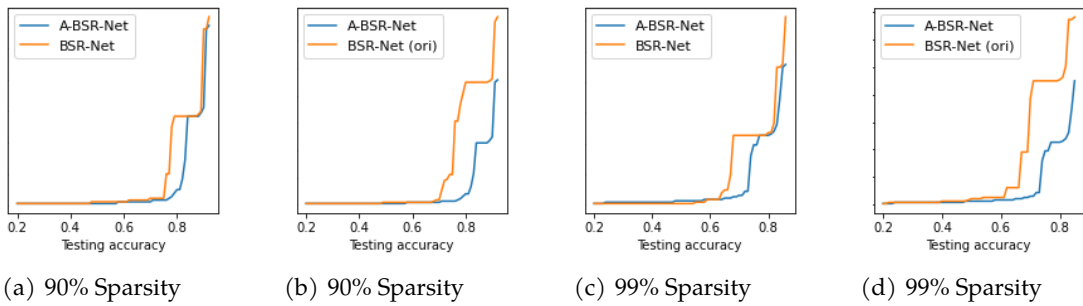


Figure 16: Comparisons (required hours to reach given accuracy. We evaluate sparse networks (99% or 90%) learned with natural training on CIFAR-10 using VGG-16.

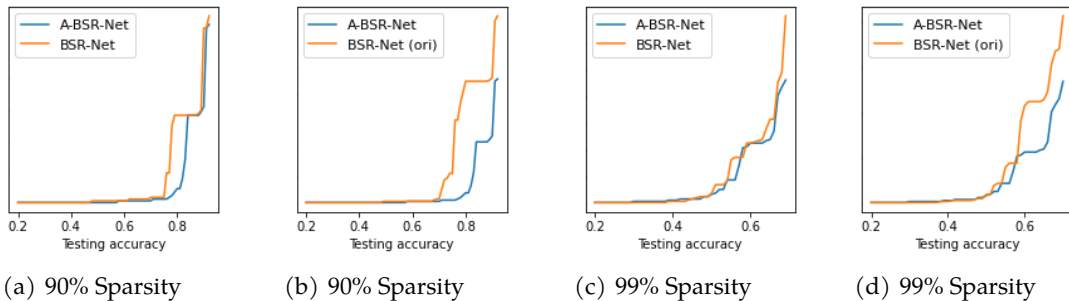


Figure 17: Comparisons (required hours to reach given accuracy). We evaluate sparse networks (99% or 90%) learned with adversarial training (objective: AT) on CIFAR-10 using VGG-16.

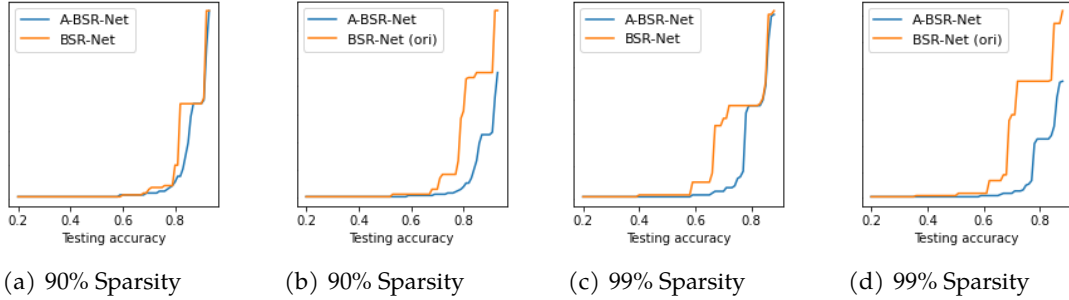


Figure 18: Comparisons (required hours to reach given accuracy). We evaluate sparse networks (99% or 90%) learned with natural training on CIFAR-10 using Wide-ResNet-28-4.

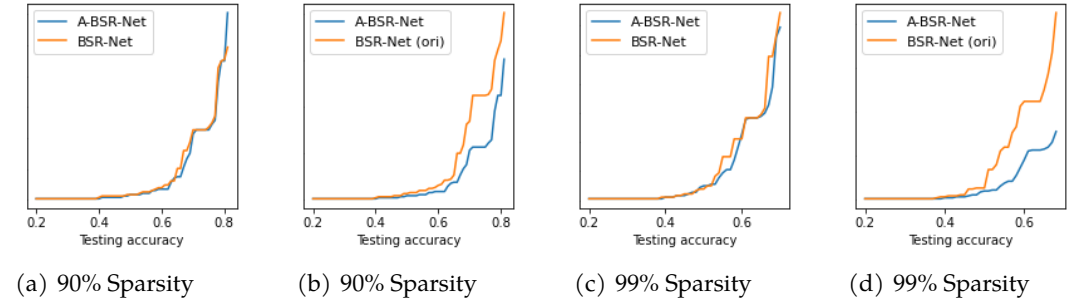


Figure 19: Comparisons (required hours to reach given accuracy). We evaluate sparse networks (99% or 90%) learned with adversarial training (objective: AT) on CIFAR-10 using Wide-ResNet-28-4.

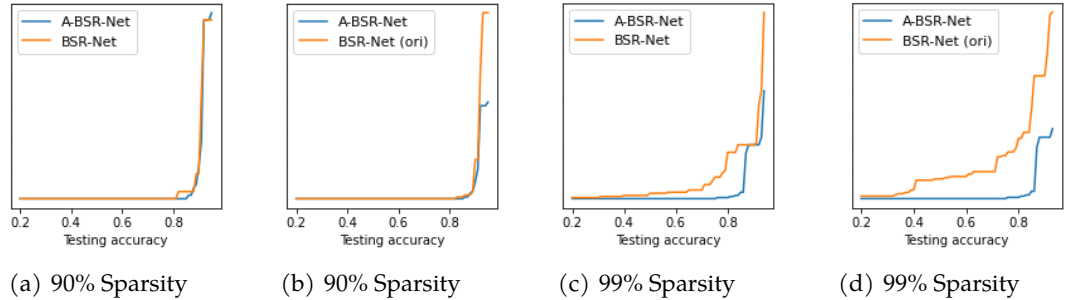


Figure 20: Comparisons (required hours to reach given accuracy). We evaluate sparse networks (99% or 90%) learned with natural training on SVHN using VGG-16.

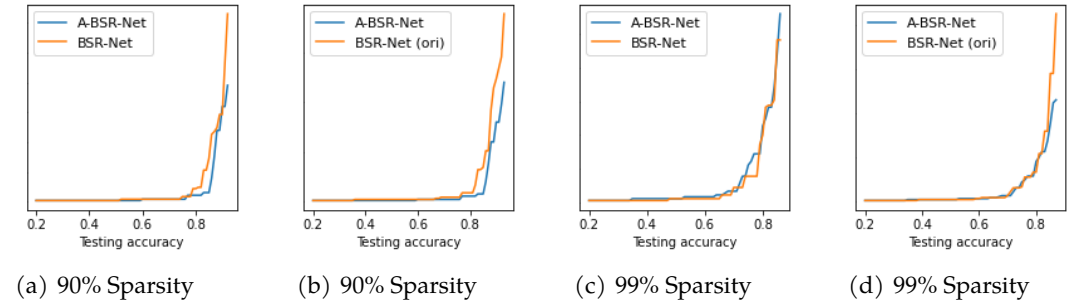


Figure 21: Comparisons (required hours to reach given accuracy). We evaluate sparse networks (99% or 90%) learned with adversarial training (objective: TRADES) on SVHN using VGG-16.

B.3. Scaling Parameter Setting

The choice of the scaling parameter γ is important to the acceleration and can be seen as a hyperparameter tuning process. We experiment with different values of γ and find that setting $\gamma = 0.1$ is a good choice for effective acceleration of training. If we tune the value of γ according to the gradient correlation of different settings, it is possible to obtain a faster convergence rate than the reported results. We first present results that are based on sparse networks (99%) learned with adversarial training (objective: AT) on CIFAR-10 using VGG-16. The sparse training method is BSR-Net.

$\gamma = 0.1$ vs $\gamma = 0.5$: As shown in Figure 22 (a), we compare the training curves (testing accuracy at different epochs) A-BSR-Net ($\gamma = 0.1$), A-BSR-Net ($\gamma = 0.5$), and BSR-Net. The yellow curve for A-BSR-Net ($\gamma = 0.5$) collapses after around 40 epochs of training, indicating a model divergence. The reason is that if setting γ close to 1, e.g., like 0.5, we will not be able to completely avoid the increase in variance. The increase in variance will lead to a decrease in performance, which is similar to "No γ " in section 5.4 of the manuscript.

$\gamma = 0.1$ vs $\gamma = 0.01$: As shown in Figure 22 (b), we compare the training curves (testing accuracy at different epochs) A-BSR-Net ($\gamma = 0.1$), A-BSR-Net ($\gamma = 0.01$), and BSR-Net. The yellow curve for A-BSR-Net ($\gamma = 0.01$) is below the blue curve for A-BSR-Net ($\gamma = 0.1$), indicating a slower convergence speed. The reason is that if γ is set small, such as 0.01, the weight of the old gradients will be small. Thus, the old gradients will have limited influence on the updated direction of the model, which tends to slow down the convergence and sometimes can lead to more training instability.

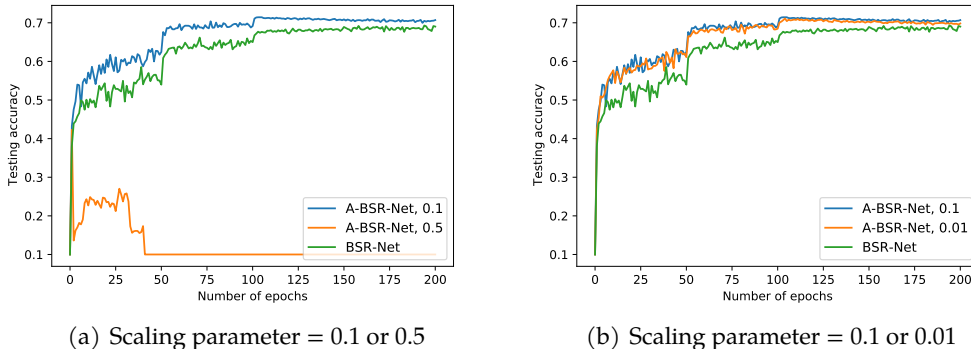


Figure 22: Comparisons (validation accuracy given the number of epochs) with different scaling parameters in BSR-Net-based models [7]. We evaluate sparse networks (99%) learned with adversarial training (objective: AT) on CIFAR-10 using VGG-16. (a) scaling parameter = 0.1 or 0.5, (b) scaling parameter = 0.1 or 0.01.

We also present results that are based on sparse networks (99%) learned with standard training on CIFAR-10 using VGG-C. The sparse training method is SET-ITOP. As shown in Figure 23, the results of setting $\gamma = 0.01, 0.5$ are similar to that of $\gamma = 0.1$, and the results of setting $\gamma = 0.9$ are worse than that of $\gamma = 0.1$. This may be due to the fact that 0.9 is too large for the relatively low gradient correlation.

B.4. Other Variance Reduction Method Comparisons

We also include more results about the comparison between our AGENT and stochastic variance reduced gradient (SVRG) [18, 19, 50], a popular variance reduction method in non-sparse case, to show the limitations of previous methods.

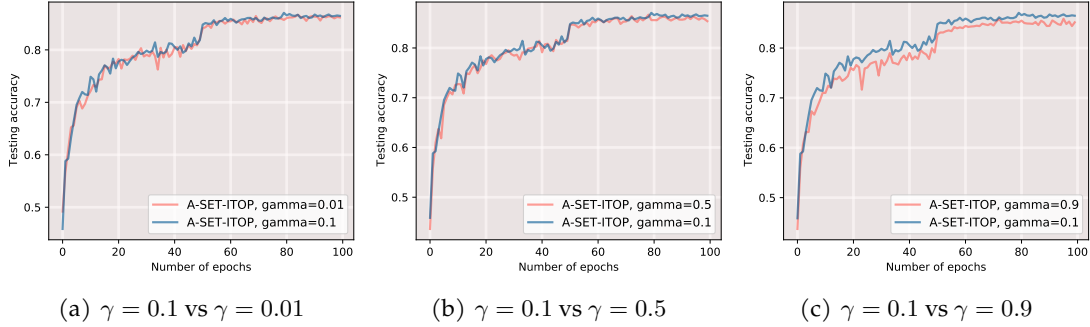


Figure 23: Comparisons (validation accuracy given the number of epochs) between different scaling factors γ . We evaluate sparse networks (99%) learned with standard training on CIFAR-10 using VGG-C where we set (a) $\gamma = 0.1$ vs $\gamma = 0.01$, (b) $\gamma = 0.1$ vs $\gamma = 0.5$, and (c) $\gamma = 0.1$ vs $\gamma = 0.9$.

B.4.1. BSR-Net-based Results

The presented results are based on sparse networks (99%) learned with adversarial training (objective: AT) on CIFAR-10 using VGG-16. As presented in Figure 24, we show the training curves (testing accuracy at different epochs) of A-BSR-Net, BSR-Net, and BSR-Net using SVRG. The yellow curve for BSR-Net using SVRG rises to around 0.4 and then rapidly decreases to a small value of around 0.1, indicating a model divergence. This demonstrates that SVRG does not work for sparse training. As for the blue curve for our A-BSR-Net, it is always above the green curve for BSR-Net, indicating a successful acceleration.

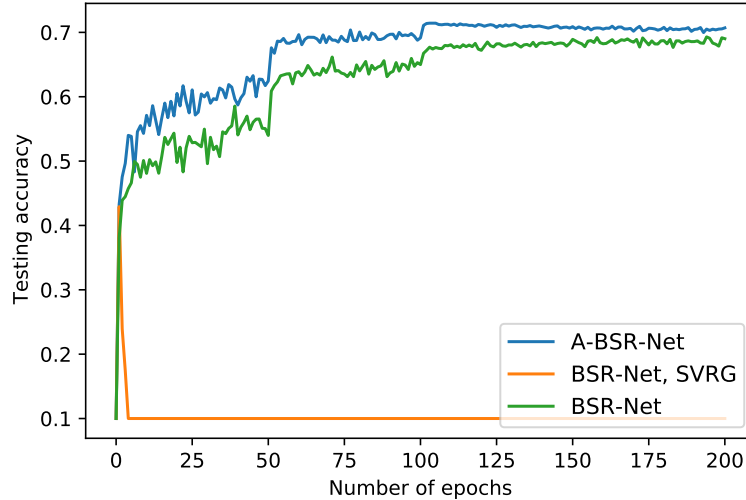


Figure 24: Comparisons (testing accuracy given the number of epochs) with different variance reduction methods in BSR-Net-based models [7]. We evaluate sparse networks (99%) learned with adversarial training (objective: AT) on CIFAR-10 using VGG-16.

B.4.2. RigL-based Results

The presented results are based on sparse networks (90%) learned with standard training on CIFAR-100 using ResNet-50. As presented in Figure 25, we show the training curves (testing accuracy at different epochs) of A-RigL, RigL, and RigL using SVRG. The yellow curve for RigL using SVRG is always below the other two curves, indicating a slower model convergence. This demonstrates that

SVRG does not work for sparse training. As for the blue curve for our A-RigL, it is always on the top of the green curve for RigL, indicating that the speedup is successful.

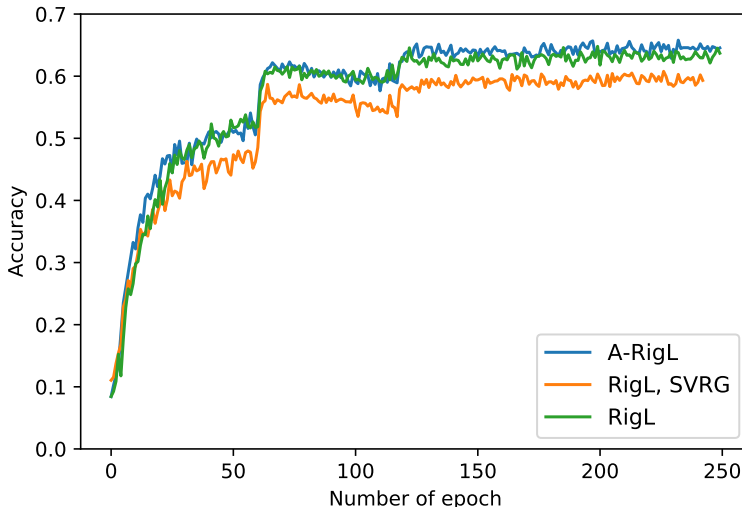


Figure 25: Comparisons (testing accuracy given the number of epochs) with different variance reduction methods in RigL-based models [2]. We evaluate sparse networks (90%) learned with standard training on CIFAR-100 using ResNet-50.

B.5. Final Accuracy Comparisons

In addition, we include more results comparing the final accuracy after sufficient training for RigL-based results on CIFAR-10/100. are shown in Table 7. Our method A-RigL tends to be the best in almost all the scenarios. This shows that our AGENT can accelerate sparse training while maintaining or even improving accuracy.

Table 7: Final accuracy (%) of RigL-based models at 0% (dense), 90% and 99% sparsity. AGENT + RigL (A-RigL) maintains or even improves the accuracy compared to that of RigL.

		DENSE	90%	99%
CIFAR-10	A-RigL	95.2 (0.24)	95.0 (0.21)	93.1 (0.25)
	RigL	95.0 (0.26)	94.2 (0.22)	92.5 (0.33)
CIFAR-100	A-RigL	72.9 (0.19)	72.1 (0.20)	66.4 (0.14)
	RigL	73.1 (0.17)	71.6 (0.26)	66.0 (0.19)

We also provide additional BSR-Net-based results for the final accuracy comparison. In addition to the BSR-Net and A-BSR-Net in the manuscript, we also include HYDRA in the appendix, which is also a SOTA sparse and adversarial training pipeline. The results are trained on SVHN using VGG-16 and WideResNet-28-4 (WRN-28-4). The final results for BSR-Net and HYDRA are obtained from [7] using their original learning rate schedules. As shown in Table 8, it is encouraging to note that our method tends to be the best in all cases when given clean test samples. In terms of robustness, our A-BSR-Net beats HYDRA in most cases, while experiencing a performance degradation compared to BSR-Net.

Table 8: Comparisons the BSR-Net [7] and HYDRA [13]. Evaluations of sparse networks learned with robust training objectives (TRADES) on SVHN using VGG-16 and WideResNet-28-4. Evaluations are after full training (200 epochs) and presented as clean/robust accuracy (%). Robust accuracy is evaluated via PGD⁵⁰ with 10 restarts $\epsilon = 8/255$.

		BSR-Net	HYDRA	Ours
90% SPARSITY	VGG-16	89.4/53.7	89.2/52.8	94.4/51.9
	WRN-28-4	92.8/55.6	94.4/43.9	95.5/46.2
99% SPARSITY	VGG-16	86.4/48.7	84.4/47.8	90.9/47.9
	WRN-28-4	89.5/52.7	88.9/39.1	92.2/51.1

B.6. Gradient Change Speed & Sparsity Level

In sparse training, when there is a small change in the weights, the gradient changes faster than in dense training, and this phenomenon can be expressed as a low correlation between the current and previous gradients, making the existing variance reduction methods ineffective.

Intuitive point of view: Considering the weights on which the current and previous gradients were calculated, there are three cases to be discussed in sparse training when the masks of current and previous gradients are different. First, if current weights are pruned, we do not need to consider their correlation because we do not need to update the current weights using the corresponding previous weights. Second, if current weights are not pruned but previous weights are pruned, the previous weights are zero and the difference between the two weights is relatively large, leading to a lower relevance. Third, if neither the current nor the previous weights are pruned, which weights are pruned can still change significantly, leading to large changes in the current and previous models. Thus, the correlation between the current and previous gradients of the weights will be relatively small. Thus, it is not a good idea to set $c = 1$ directly in sparse training which can even increase the variance and slow down the convergence.

When the masks of the current and previous gradients are the same, the correlation still tends to be weaker. As we know, $c_t^* = \frac{\text{Cov}(g_{\text{new}}, g_{\text{old}})}{\text{Var}(g_{\text{old}})}$. Even if $\text{Cov}(g_{\text{new}}, g_{\text{old}})$ does not decrease, the variance $\text{Var}(g_{\text{old}})$ increases in sparse training, leading to a decrease in c_t^* .

Apart from the analysis above, we also do some experiments to demonstrate that the gradient changes faster as the sparsity increases. To measure the rate of change, our experiments are described below.

Correlation over the course of training: We also analyze the gradient correlation during the standard training of sparse ResNet-50 on CIFAR-100 using RigL. The results are summarized in Figure 26, where the blue curves represent the gradient correlation of dense training (0% sparsity) and the pink curves denote the correlation of sparse training, i.e., RigL. As we can see, the correlation between dense and sparse training is close. For 80% sparsity, sparse training tends to have a lower correlation compared to dense training, especially in late training stages. For 90% and 95%, sparse training also gives lower relevance than dense training, and the differences become larger with increasing sparsity.

Correlation of the fully-trained model: We begin with fully-trained checkpoints from ResNet-50 on CIFAR-100 with RigL and SET at 0%, 50%, 80%, 90%, and 95% sparsity. We calculate and store the gradient of each weight on all training data. Then, we add Gaussian perturbations (std = 0.015) to all the weights and calculate the gradients again. Lastly, we calculate the correlation between the gradient of the new perturbed weights and the old original weights.

As we know, there is always a difference between the old and new weights. If the gradients become very different after adding some small noise to the weights, the new and old gradients will tend to have smaller correlations. If the gradients do not change a lot after adding some small noise, the



Figure 26: Gradient correlation in dense training and SET-ITOP. We evaluate sparse networks learned with standard training on CIFAR-10 using VGG-C. We compare the correlation between dense training (sparsity=0%) and sparse training in sparsity (a) 50%, (b) 80%, (c) 90%, and (d) 95%.

old and new gradients will have a higher correlation. Thus, we add Gaussian noise to the weights to simulate the difference between the new and old gradients. As shown in Table 9, the correlation decreases with increasing sparsity, which indicates a weaker correlation in sparse training and supports our claim.

Table 9: Correlation between the gradient of the new perturbed weights and the old original weights from ResNet-50 on CIFAR-100 produced by RigL and SET at different sparsity including 0%, 50%, 80%, 90%, 95%, 99%.

SPARSITY	0%	50%	80%	90%	95%
RESNET-50, CIFAR-100 (RigL)	0.6005	0.4564	0.3217	0.1886	0.1590
RESNET-50, CIFAR-100 (SET)	0.6005	0.4535	0.2528	0.1763	0.1195

B.7. Variants of RigL

RigL is one of the most popular dynamic sparse training pipelines which uses weight magnitude for pruning and gradient magnitude for growing. Our method adaptively updates the new batch gradient using the old storage gradient which usually has less noise. As a result, the variance of the new batch gradient is reduced, leading to fast convergence. Currently, we only use gradients with corrected variance in weight updates. A natural question is how it performs if we also use this variance-corrected gradient for weight growth in RigL.

We do some experiments in RigL-based models trained on CIFAR-10. As shown in Figure 27, the blue curves (RigL-ITOP-G) and yellow curves (RigL-ITOP) correspond to the weight growth with and without the variance-corrected gradient, respectively. We can see that in the initial stage, the blue curves are higher than the yellow curves. But after the first learning rate decay, they tend to be lower than the yellow curves. This suggests that weight growth using a variance-corrected gradient

at the beginning of training can help the model improve accuracy faster. However, this may lead to a slight decrease in accuracy in the later training stages. This may be due to the fact that some variance in the gradient can help the model explore local regions better and find better masks as the model approaches its optimal point.

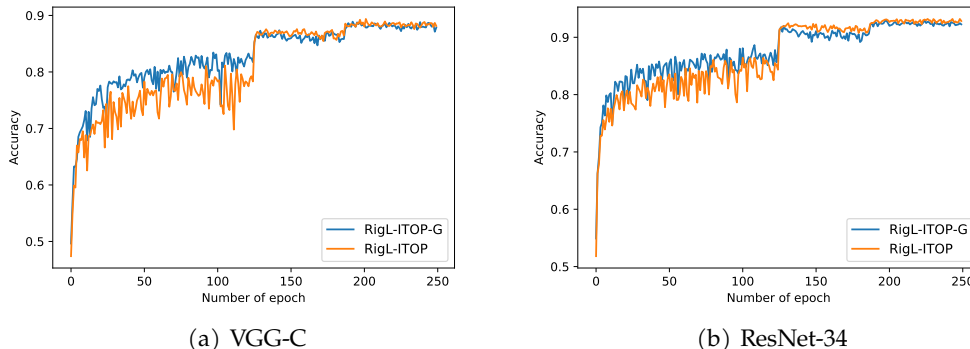


Figure 27: Comparisons (testing accuracy given the number of epochs) between weight growth with (RigL-ITOP-G) and without (RigL-ITOP) variance-corrected gradient [6]. We evaluate sparse networks (99%) learned with standard training on CIFAR-10 using (a) VGG-C and (b) ResNet-34.

B.8. Comparison with Reducing Learning Rate

To demonstrate the design of the scaling parameter γ , we compare our AGENT with "Reduce LR", where we remove the scaling parameter γ from AGENT and set the learning rate to 0.1 times the original one. As shown in Table 10, reducing the learning rate can lead to a comparable convergence rate in the early stage. However, it slows down the later stages of training and leads to sub-optimal final accuracy. The reason is that it reduces both signal and noise and therefore does not improve the signal-to-noise ratio or speed up the sparse training.

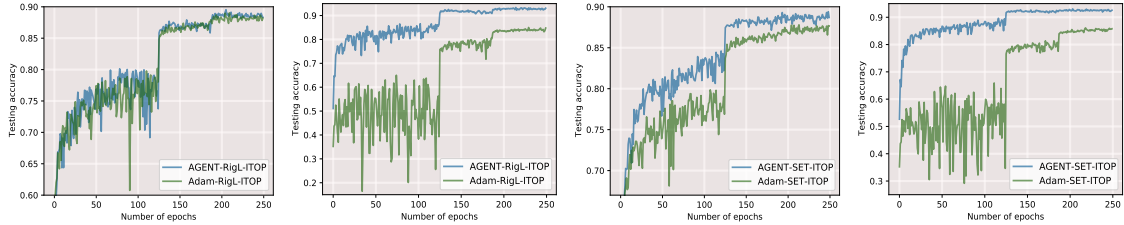
The motivation of γ is to avoid introducing large variance due to error in approximating c_t and bias due to the adversarial training. The true correlation depends on many factors such as the dataset, architecture, and sparsity. In some cases, it can be greater or smaller than 10%. For the value of γ , it is a hyperparameter and we can choose different values for different settings. In our case, for simplicity, we choose $\gamma = 0.1$ for all the settings, and find that it works well and accelerates the convergence. If we tune the value of γ for different settings according to their corresponding correlations, it is possible to obtain faster convergence rates.

Table 10: Testing accuracy (%) of SET-ITOP-based models for AGENT (ours) and "Reduce LR". Sparse VGG-C and ResNet-34 are learned in standard setups.

EPOCH	20	80	130	180	240
REDUCE LR (VGG-C, SET-ITOP)	76.5	81.3	84.6	85.5	85.5
AGENT (VGG-C, SET-ITOP)	76.1	81.5	87.6	87.1	88.6
REDUCE LR (RESNET-34, SET-ITOP)	81.4	85.9	89.3	89.5	89.8
AGENT (RESNET-34, SET-ITOP)	83.0	85.6	92.0	92.3	92.5

B.9. Comparison with Momentum-based Methods

To some extent, our AGENT is designed with a similar idea to the momentum-based method [53, 54], where old gradients are used to improve the current batch gradient. The momentum-based approach works well in dense settings. However, the momentum-based method still suffers from optimization difficulties due to sparsity constraints. The reason is that it does not take into account sparse and adversarial training characteristics such as the reduced correlation between current and



(a) VGG-C(RigL) (b) ResNet-34(RigL) (c) VGG-C(SET) (d) ResNet-34(SET)

Figure 28: Testing accuracy for ITOP-based models at 99% sparsity on CIFAR-10. AGENT-based training (blue curves) converges faster than Adam-based training (pink curves).

previous gradients and potential bias of gradient estimator, and fails to provide an adaptive balance between old and new information. When the correlation is low, the momentum-based method can still incorporate too much of the old information and increase the gradient variance or bias. In contrast, our AGENT is designed for sparse and adversarial training and can establish finer adaptive control over how much information we should take from the old to help the new.

For example, in our baseline SGD, following the original code base, we have also added momentum to the optimizer. However, as shown in the pink curves in Figure 2, it still has training instability and convergence problems. The reason is that they do not take into account the sparse and adversarial training characteristics and cannot provide an adaptive balance between old and new information.

Our method AGENT is designed for sparse and adversarial training and can establish a finer control over how much information we should get from the old to help the new. To demonstrate the importance of this fine-grained adaptive balance, we do ablation studies in Section 6.4. In "Fixed c_t ", we set $c_t = 0.1$ and test the convergence rate without the adaptive control. We find that the adaptive balance (ours) outperforms "Fixed c_t " in almost all cases, especially in adversarial training. For standard training, "Fix c_t " provides similar convergence rates to our method, while ours tends to have better final scores.

B.10. Comparison with Other Adaptive Gradient Methods

We also compare our AGENT with other adaptive gradient methods, where we take Adam [55] as an example. As shown in Figure 28, AGENT-RigL-ITOP and AGENT-SET-ITOP (blue curves) are usually above Adam-RigL-ITOP and Adam-SET-ITOP (green curves), indicating that our AGENT converges faster compared to Adam. This demonstrates the importance of using correlation in sparse training to balance old and new information.

B.11. Different Total Number of Training Epochs

In this section, we show that our method can achieve acceleration over different training budgets (i.e., number of training epochs), rather than being a pseudo-proposition of better early performance compared to the baseline method. To demonstrate this, we add experiments under different total number of training epochs and change the learning rate scheduler accordingly to allow convergence.

Take the SET-ITOP as an example. In the main paper, we follow the baseline paper where the epoch number is 250 and the learning rate scheduler is set as the stepwise learning rate with decay points 125 (i.e., 0.5×250) and 187 (i.e., 0.75×250). To reduce the epoch number and allow convergence, we set the epoch number as 50 and 100 where the decay points are set as $\{25, 37\}$ and $\{50, 75\}$, respectively. As shown in Figure 29, blue curves (our A-SET-ITOP) are usually on top of pink curves (SET-ITOP), implying acceleration from our AGENT.

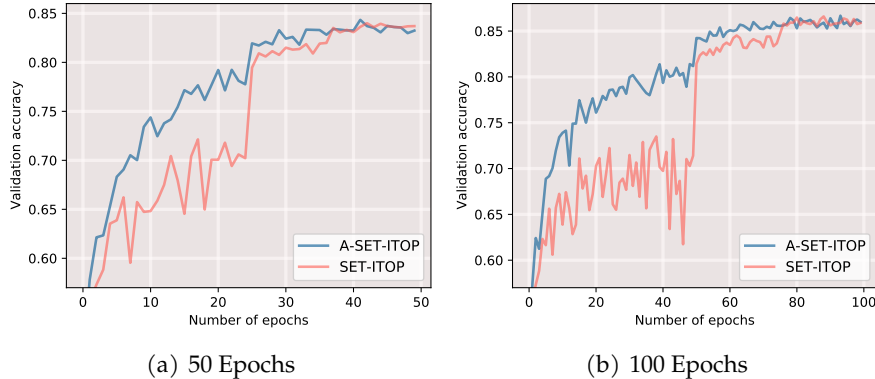


Figure 29: Comparisons (validation accuracy given the number of epochs) between A-SET-ITOP and SET-ITOP. We evaluate sparse networks (99%) learned with standard training on CIFAR-10 using VGG-C under (a) 50 training epochs, and (b) 100 training epochs.

B.12. Smoothing Factor Tuning

For smoothing factor α , we follow the default value in Deng et al. [52] which is set as 0.3. We add some experiments to test the influence of α . We further compare the validation accuracy across different smoothing factors α . As shown in Figure 30, the results of setting $\alpha = 0.05, 0.5, 0.9$ are similar to that of $\alpha = 0.3$. Thus, our method is not sensitive to the choice of α , and we can follow the default 0.3.

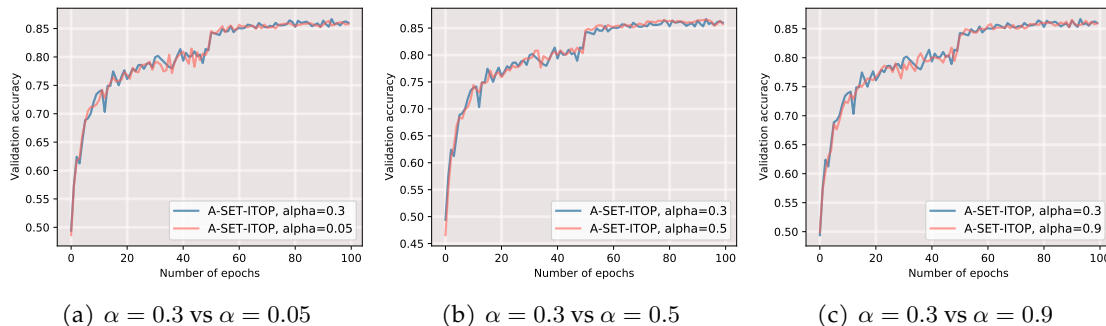


Figure 30: Comparisons (validation accuracy given the number of epochs) between different smoothing factors α . We evaluate sparse networks (99%) learned with standard training on CIFAR-10 using VGG-C where we set (a) $\alpha = 0.3$ vs $\alpha = 0.05$, (b) $\alpha = 0.3$ vs $\alpha = 0.5$, and (c) $\alpha = 0.3$ vs $\alpha = 0.9$.

B.13. Fixed c_t Tuning

We add a more realistic baseline of "Fixed c_t " with good hyperparameter tuning to show that adaptive re-weighting is crucial. The term "Fixed c_t " corresponds to fixing weight $c_t = 0.1$ during training, which is mentioned in our ablation studies in Section 6.6. Specifically, we further check different c_t in "Fixed c_t " and compare their validation accuracy with our A-SET-ITOP. As shown in Figure 31, when c_t is fixed as 0.001, 0.001, and 0.1, the pink curves ("Fixed c_t ") are lower than the blue curve (A-SET-ITOP) in the early stages, indicating slower early convergence in "Fixed c_t ". When fixing c_t as 0.5, 0.8, and 1.0, the whole pink curves ("Fixed c_t ") are below the blue curves (A-SET-ITOP), implying slower convergence in "Fixed c_t ".

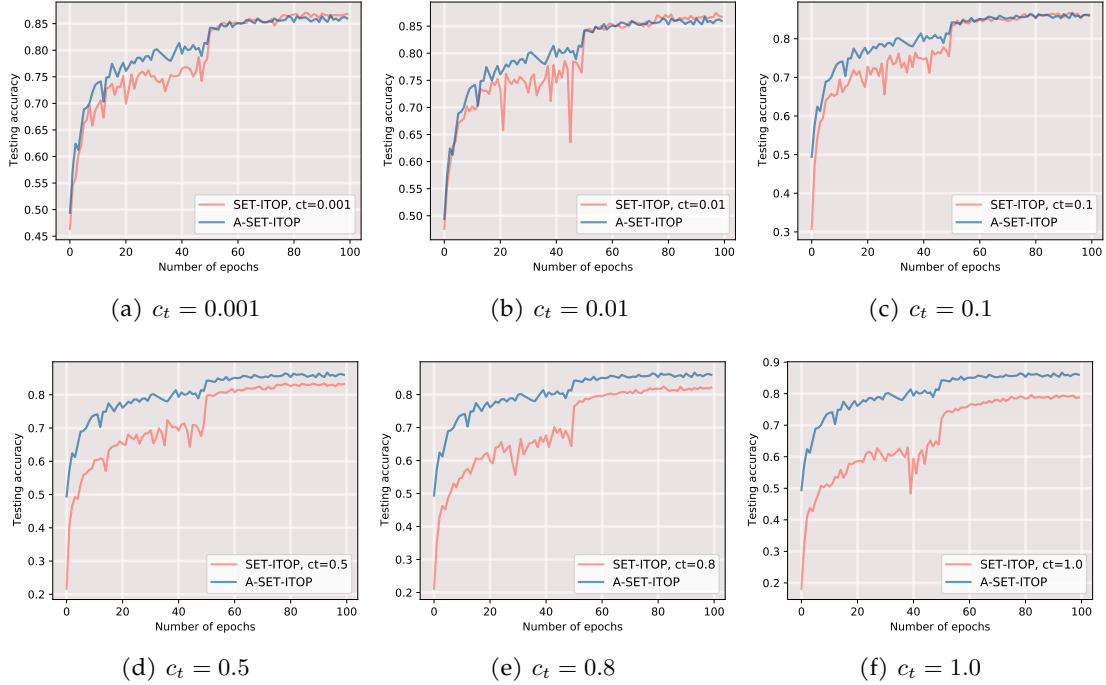


Figure 31: Comparisons (validation accuracy given the number of epochs) between different c_t in "Fixed c_t ". We evaluate sparse networks (99%) learned with standard training on CIFAR-10 using VGG-C where we compare A-SET-ITOP with (a) $c_t = 0.001$, (b) $c_t = 0.01$, (c) $c_t = 0.1$, (d) $c_t = 0.5$, (e) $c_t = 0.8$, (f) $c_t = 1.0$

B.14. Loss Value Comparisons

Apart from accuracy, we also include loss comparison to demonstrate the acceleration. As shown in Figure 32, the blue curves for our A-SET-ITOP are usually below the pink curves for SET-ITOP, implying successful acceleration.

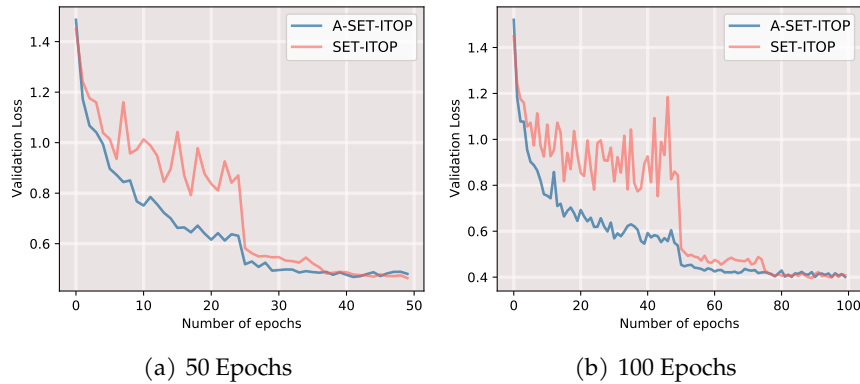


Figure 32: Comparisons (validation loss given the number of epochs) between A-SET-ITOP and SET-ITOP. We evaluate sparse networks (99%) learned with standard training on CIFAR-10 using VGG-C under (a) 50 training epochs, and (b) 100 training epochs.

B.15. More Baseline Comparison

We add more results where ADAM and SVGR are compared together. As shown in Figure 33, the blue curves, pink curves, and green curves represent our AGENT, Adam, and SVRG, respectively. The blue curves of our AGENT are usually higher than the pink and green curves, indicating faster convergence using our AGENT compared to the other two methods.

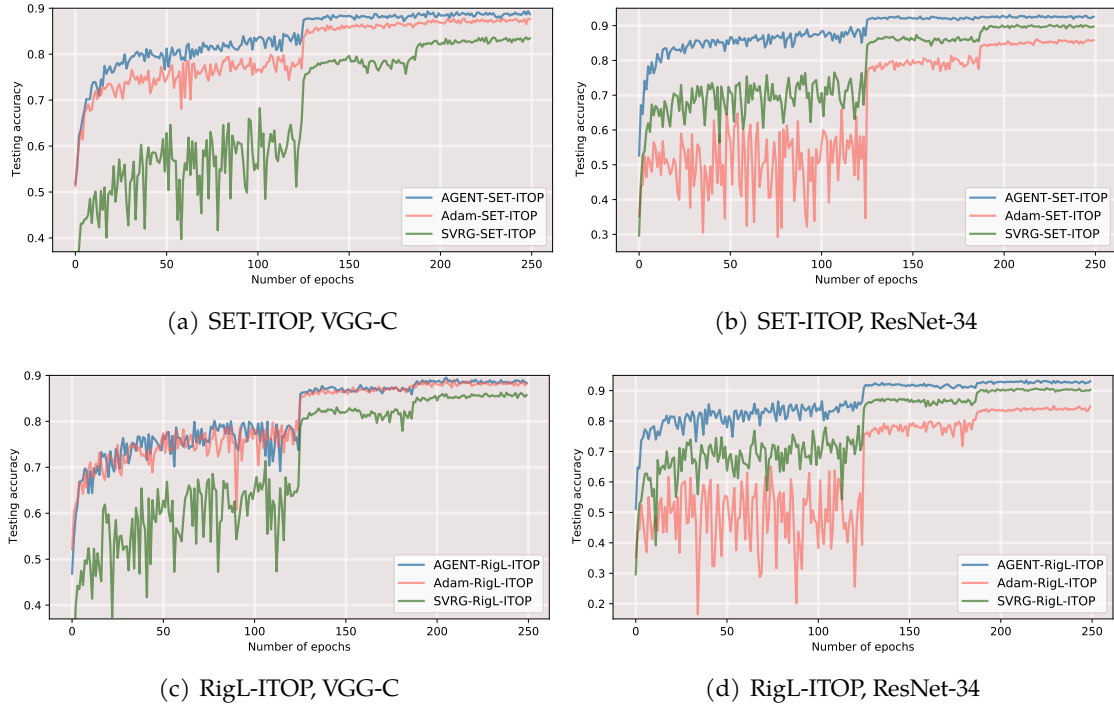


Figure 33: Comparison between our AGENT, Adam, and SVRG. We evaluate sparse networks learned with standard training on CIFAR-10. (a) SET-ITOP, VGG-C, (b) SET-ITOP, ResNet-34, (c) RigL-ITOP, VGG-C, and (d) RigL-ITOP, ResNet-34.

B.16. More Training Time Comparison

We check the training time of our method and baseline methods. For ITOP-based results, the training time ratio between our A-SET-ITOP and SET-ITOP is 5:3, and the training time ratio between our A-RigL-ITOP and RigL-ITOP is 2:1. For BSR-Net based results, the training time ratio between our A-BSR-Net and BSR-Net is 5:4. Despite our current training time does not have advantages over baseline methods, our training time can be easily reduced by the following ways.

- We can use sparse gradients in sparse training, which effectively reduces the cost of backward in sparse training and can be easily applied to our method [51].
- We can use parallel computing. Since the additional forward and backward over the old model parameters are fully parallelizable, we can view it as doubling the mini-batch size [49].
- We can follow the idea of SAGA and store gradients for each sample. Then, we do not need extra forward and backward steps, saving the wall-clock time [40].

B.17. Gradient Norm Comparison

Larger gradient norms are important for sparse training [65]. We conduct experiments and find that our AGENT can improve the gradient norm. Specifically, we train 99% sparse VGG-C and ResNet-34 on CIFAR-10. We compare AGENT + RigL-ITOP (A-RigL-ITOP) and RigL-ITOP, as well as AGENT + SET-ITOP (A-SET-ITOP) and SET-ITOP. For the gradient norm, we calculate the average gradient norm for each training phase, i.e., 1st to 50th epochs, 51st to 100th epochs, 101st to 150th epochs, 151st to 200th epochs, and 201st to 250th epochs. The results are summarized in the table below. Our AGENT can slightly improve the gradient norm, which is important for sparse training.

Table 11: Gradient norm of RigL-ITOP-based and SET-ITOP-based models for AGENT (ours) and sparse training baseline methods. 99% Sparse VGG-C is learned in standard setups.

	1ST TO 50TH	51ST TO 100TH	101ST TO 150TH	151ST TO 200TH	201ST TO 250TH
RigL-ITOP	3.23	2.40	2.19	3.04	3.47
A-RigL-ITOP	3.19	2.43	2.29	3.16	3.54
SET-ITOP	3.25	2.79	2.89	4.21	4.62
A-SET-ITOP	3.27	2.85	2.96	4.22	4.88

C. Additional Details about Experiment Settings

C.1. Gradient Variance and Correlation Calculation

We calculate the gradient variance and correlation of the ResNet-50 on CIFAR-100 from RigL [2] and SET [1] at different sparsities including 0%, 50%, 80%, 90%, and 95%. The calculation is based on the checkpoints from Sundar and Dwaraknath [66].

Gradient variance: We first load fully trained checkpoints for the 0%, 50%, 80%, 90%, and 95% sparse models. Then, to see the gradient variance around the converged optimum, we add small perturbations to the weights and compute the mean of the gradient variance. For each checkpoint, we do three replicates.

Gradient correlation: We begin with fully-trained checkpoints at 0%, 50%, 80%, 90%, and 95% sparsity. We calculate and store the gradient of each weight on all training data. Then, we add Gaussian perturbations to all the weights and calculate the gradients again. Lastly, we calculate the correlation between the gradient of the new perturbed weights and the old original weights. For each checkpoint, we do three replicates.

C.2. Implementations

In BSR-Net-based results, aligned with the choice of Özdenizci and Legenstein [7], the gradients for all models are calculated by SGD with momentum and decoupled weight decay [67]. All models are trained for 200 epochs with a batch size of 128.

In RigL-based results, we follow the settings in Evci et al. [2], Sundar and Dwaraknath [66]. We train all the models for 250 epochs with a batch size of 128, and parameters are optimized by SGD with momentum.

In ITOP-based results, we follow the settings in Liu et al. [6]. For CIFAR-10 and CIFAR-100, we train all the models for 250 epochs with a batch size of 128. For ImageNet-2012, we train all the models for 100 epochs with a batch size of 64. Parameters are optimized by SGD with momentum.

C.3. Learning Rate

Aligned with popular sparse training methods [2, 6, 7], we choose piecewise constant decay schedulers for learning rate and weight decay. In our A-BSR-Net, we use the 50th and 100th epochs as the dividing points of our learning rate decay scheduler. The reason is that our approach has faster convergence and doesn't require a long warm-up period. In the evaluation shown in the manuscript, we also use this scheduler for BSR-Net for a more accurate and fair comparison.

C.4. Initialization (BSR-Net-based results)

Consistent with [7], we also choose Kaiming initialization to initialize the network weights [68]

C.5. Benchmark Datasets (BSR-Net-based results)

For a fair comparison, we choose the same benchmark datasets as [7]. Specifically, we use CIFAR-10 and CIFAR-100 [57] and SVHN [58] in our experiments. Both CIFAR-10 and CIFAR-100 datasets include 50,000 training and 10,000 test images. SVHN dataset includes 73,257 training and 26,032 test samples.

C.6. Data Augmentation

We follow a popular data augmentation method used in [7, 61]. In particular, we randomly shift the images to the left or right, crop them back to their original size, and flip them in the horizontal direction. In addition, all the pixel values are normalized in the range of $[0, 1]$.

D. Sparse Training Method Description

D.1. Sparse Training Overview

As the cost of deep neural networks (DNNs) increases, there is a growing interest in efficiency issues, such as model efficiency and data efficiency [69–73]. To improve model efficiency, sparsity is a common class of solutions, and it has been found that there is a sparsity pattern in the parameters of trained DNNs [74–76]. Sparse training is a popular method for introducing sparsity into DNNs and achieving resource efficiency in DNNs. Specifically, to obtain a 90% sparse DNN, we randomly initialize a 90% sparse DNN. Then, we maintain sparse weights throughout the training process, pruning and regrowing a certain number of weights every m iterations. Thus, we can save training memory and produce sparse models with dense performance levels. Several widely used sparse training methods are described below.

D.2. Sparse Training Method: SET

SET [1] is a broadly-used sparse training method that prunes and regrows connections by examining the magnitude of the weights.

D.3. Sparse Training Method: RigL

RigL [2] is another popular dynamic sparse training method that uses weight and gradient magnitudes to learn the connections.

D.4. Sparse Training Method: BSR-Net

Bayesian Sparse Robust Training (BSR-Net) [7] is a Bayesian Sparse and Robust training pipeline. Based on a Bayesian posterior sampling principle, a network rewiring process simultaneously learns the sparse connectivity structure and the robustness-accuracy trade-off based on the adversarial learning objective. More specifically, regarding its mask update, it prunes all negative weights and grows new weights randomly.

D.5. Sparse Training Method: ITOP

ITOP [6] is another recent pipeline for dynamic sparse training, which uses sufficient and reliable parameter exploration to achieve in-time over-parameterization and find well-performing sparse models.

E. Limitations of Our Adaptive Gradient Correction Method

E.1. Extra FLOPs

Similar to SVRG, our AGENT increases the training FLOPs in each iteration due to the extra forward and backward used to compute the old gradients.

However, the true computation difference can be smaller and the GPU-based running time of SVRG will not be affected that much. For example, in the adversarial setting, we need additional computations to generate the adversarial samples, which is time-consuming and only needs to be done once in each iteration of our AVR and SGD. For BSR-Net, we empirically find that the ratio of time required for each iteration of our AVR and SGD is about 1.2.

There are also several methods to reduce the extra computation caused by SVRG. The first approach is to use the sparse gradients proposed by M Elibol (2020) [51]. It can effectively reduce the computational cost of SVRG and can be easily applied to our method. The second approach is suggested by Allen-Zhu and Hazan (2016) [49]. The extra cost of computing batch gradient on old model parameters is totally parallelizable. Thus, we can view SVRG as doubling the mini-batch size. Third, we can follow the idea of SAGA [40] and store gradients for individual samples. By this way, we do not need the extra forward and backward step and save the computation. But it requires extra memory to store the gradients.

In the main manuscript, we choose to compare the convergence speed of our ADSVRG and SGD for the same number of pass data (epoch), which is widely used as a criterion to compare SVRG-based optimization and SGD [18, 22, 46, 49]. A comparison in this way can demonstrate the accelerating effect of the optimization method and provide inspiration for future work.

E.2. Inefficiencies in Large Data or Models

When the amount of data is large, the additional gradient computation in our AGENT can be very extensive. But we can randomly sample a small portion of the data and compute \tilde{g} on that small portion, thus saving the extra computation significantly. Specifically, we can sample a subdata that is larger than the batch data, but smaller than the full data. When we compute \tilde{g} on this small portion of data, there will be less noise compared to computing the gradient on the batch data. Therefore, we can use \tilde{g} on this small portion in AGENT to correct the gradient and accelerate sparse training. Take ImageNet-2012 as a large data example. We use a batch size of 128. In each epoch, we compute \tilde{g} on sub-data of size 50,000, instead of computing \tilde{g} on the full data of size 1.28 million.

When the number of model parameters is large, the extra memory burden of storing and accessing additional gradients can be significant. But we can significantly save the memory via a sparse matrix [77]. During sparse training, we only need to update the active weights. Thus, we only need to store the gradients for the active weights. The stored gradient is sparse and can be stored in the sparse matrix to save memory.

E.3. Scaling Parameter Tuning

In our adaptive variance reduction method (AVR), we add an additional scaling parameter γ which needs to be adjusted. We find that setting $\gamma = 0.1$ is a good choice for BSR-Net, RigL, and ITOP. However, it can be different for other different sparse training pipelines.

E.4. Robust Accuracy Degradation

For the final accuracy results of BSR-Net-based models, there is a small decrease in the robustness accuracy after using our AVR. It is still an open question of how to further improve the robust accuracy when using adaptive variance reduction in sparse and adversarial training.